

```
1: =====
2: F.A.R.F.A.N PIPELINE CODE AUDIT - BATCH 20
3: =====
4: Generated: 2025-12-07T06:17:25.812595
5: Files in this batch: 17
6: =====
7:
8:
9: =====
10: FILE: src/farfan_pipeline/processing/aggregation.py
11: =====
12:
13: """
14: Aggregation Module - Hierarchical Score Aggregation System
15:
16: This module implements the complete aggregation pipeline for the policy analysis system:
17: - FASE 4: Dimension aggregation (60 dimensions: 6 → 27 10 policy areas)
18: - FASE 5: Policy area aggregation (10 areas)
19: - FASE 6: Cluster aggregation (4 MESO questions)
20: - FASE 7: Macro evaluation (1 holistic question)
21:
22: Requirements:
23: - Validation of weights, thresholds, and hermeticity
24: - Comprehensive logging and abortability at each level
25: - No strategic simplification
26: - Full alignment with monolith specifications
27: - Uses canonical notation for dimension and policy area validation
28:
29: Architecture:
30: - DimensionAggregator: Aggregates 5 micro questions → 1 dimension score
31: - AreaPolicyAggregator: Aggregates 6 dimension scores → 1 area score
32: - ClusterAggregator: Aggregates multiple area scores → 1 cluster score
33: - MacroAggregator: Aggregates all cluster scores → 1 holistic evaluation
34: """
35:
36: from __future__ import annotations
37:
38: import logging
39: from collections import defaultdict
40: from dataclasses import dataclass, field
41: from typing import TYPE_CHECKING, Any, TypeVar
42: from farfan_pipeline.core.calibration.parameter_loader import get_parameter_loader
43: from farfan_pipeline.core.calibration.decorators import calibrated_method
44: from farfan_pipeline.core.parameters import ParameterLoaderV2
45:
46: # SOTA imports
47: from farfan_pipeline.processing.aggregation_provenance import (
48:     AggregationDAG,
49:     ProvenanceNode,
50: )
51: from farfan_pipeline.processing.uncertainty_quantification import (
52:     BootstrapAggregator,
53:     UncertaintyMetrics,
54:     aggregate_with_uncertainty,
55: )
56: from farfan_pipeline.processing.choquet_adapter import (
```

```
57:     ChoquetProcessingAdapter,
58:     create_default_choquet_adapter,
59: )
60:
61: if TYPE_CHECKING:
62:     from collections.abc import Callable, Iterable
63:
64: T = TypeVar('T')
65:
66:
67: @dataclass(frozen=True)
68: class AggregationSettings:
69:     """Resolved aggregation settings derived from the questionnaire monolith."""
70:
71:     dimension_group_by_keys: list[str]
72:     area_group_by_keys: list[str]
73:     cluster_group_by_keys: list[str]
74:     dimension_question_weights: dict[str, dict[str, float]]
75:     policy_area_dimension_weights: dict[str, dict[str, float]]
76:     cluster_policy_area_weights: dict[str, dict[str, float]]
77:     macro_cluster_weights: dict[str, float]
78:     dimension_expected_counts: dict[tuple[str, str], int]
79:     area_expected_dimension_counts: dict[str, int]
80:
81:     @classmethod
82:     def from_monolith(cls, monolith: dict[str, Any] | None) -> AggregationSettings:
83:         """Build aggregation settings from canonical questionnaire data."""
84:         if not monolith:
85:             return cls(
86:                 dimension_group_by_keys=["policy_area", "dimension"],
87:                 area_group_by_keys=["area_id"],
88:                 cluster_group_by_keys=["cluster_id"],
89:                 dimension_question_weights={},
90:                 policy_area_dimension_weights={},
91:                 cluster_policy_area_weights={},
92:                 macro_cluster_weights={},
93:                 dimension_expected_counts={},
94:                 area_expected_dimension_counts={}
95:             )
96:
97:         blocks = monolith.get("blocks", {})
98:         niveles = blocks.get("niveles_abstraccion", {})
99:         policy_areas = niveles.get("policy_areas", [])
100:        clusters = niveles.get("clusters", [])
101:        micro_questions = blocks.get("micro_questions", [])
102:
103:        aggregation_block = (
104:            monolith.get("aggregation")
105:            or blocks.get("aggregation")
106:            or monolith.get("rubric", {}).get("aggregation")
107:            or {}
108:        )
109:
110:        # Map question_id \206\222 base_slot for later normalization
111:        question_slot_lookup: dict[str, str] = {}
112:        dimension_slot_map: dict[str, set[str]] = defaultdict(set)
```

```
113:     dimension_expected_counts: dict[tuple[str, str], int] = defaultdict(int)
114:
115:     for question in micro_questions:
116:         qid = question.get("question_id")
117:         dim_id = question.get("dimension_id") or question.get("dimension")
118:         area_id = question.get("policy_area_id") or question.get("policy_area")
119:         base_slot = question.get("base_slot")
120:
121:         if dim_id and qid and not base_slot:
122:             base_slot = f"{dim_id}-{qid}"
123:
124:         if qid and base_slot:
125:             question_slot_lookup[qid] = base_slot
126:             dimension_slot_map[dim_id].add(base_slot)
127:
128:         if area_id and dim_id:
129:             dimension_expected_counts[(area_id, dim_id)] += 1
130:
131:     area_expected_dimension_counts: dict[str, int] = {}
132:     for area in policy_areas:
133:         area_id = area.get("policy_area_id") or area.get("id")
134:         if not area_id:
135:             continue
136:         dims = area.get("dimension_ids") or []
137:         area_expected_dimension_counts[area_id] = len(dims)
138:
139:     group_by_block = aggregation_block.get("group_by_keys") or {}
140:     dimension_group_by_keys = cls._coerce_str_list(
141:         group_by_block.get("dimension"),
142:         fallback=["policy_area", "dimension"],
143:     )
144:     area_group_by_keys = cls._coerce_str_list(
145:         group_by_block.get("area"),
146:         fallback=["area_id"],
147:     )
148:     cluster_group_by_keys = cls._coerce_str_list(
149:         group_by_block.get("cluster"),
150:         fallback=["cluster_id"],
151:     )
152:
153:     dimension_question_weights = cls._build_dimension_weights(
154:         aggregation_block.get("dimension_question_weights") or {},
155:         question_slot_lookup,
156:         dimension_slot_map,
157:     )
158:     policy_area_dimension_weights = cls._build_area_dimension_weights(
159:         aggregation_block.get("policy_area_dimension_weights") or {},
160:         policy_areas,
161:     )
162:     cluster_policy_area_weights = cls._build_cluster_weights(
163:         aggregation_block.get("cluster_policy_area_weights") or {},
164:         clusters,
165:     )
166:     macro_cluster_weights = cls._build_macro_weights(
167:         aggregation_block.get("macro_cluster_weights") or {},
168:         clusters,
```

```
169:         )
170:
171:     return cls(
172:         dimension_group_by_keys=dimension_group_by_keys,
173:         area_group_by_keys=area_group_by_keys,
174:         cluster_group_by_keys=cluster_group_by_keys,
175:         dimension_question_weights=dimension_question_weights,
176:         policy_area_dimension_weights=policy_area_dimension_weights,
177:         cluster_policy_area_weights=cluster_policy_area_weights,
178:         macro_cluster_weights=macro_cluster_weights,
179:         dimension_expected_counts=dict(dimension_expected_counts),
180:         area_expected_dimension_counts=area_expected_dimension_counts,
181:     )
182:
183:     @staticmethod
184:     def _coerce_str_list(value: Any, *, fallback: list[str]) -> list[str]:
185:         if isinstance(value, list) and all(isinstance(item, str) for item in value):
186:             return value or fallback
187:         return fallback
188:
189:     @staticmethod
190:     def _normalize_weights(weight_map: dict[str, float]) -> dict[str, float]:
191:         if not weight_map:
192:             return {}
193:         # Discard negative weights and normalize remaining ones
194:         positive_map = {k: float(v) for k, v in weight_map.items() if isinstance(v, (float, int)) and float(v) >= 0.0}
195:         if not positive_map:
196:             equal = 1.0 / len(weight_map)
197:             return {k: equal for k in weight_map}
198:         total = sum(positive_map.values())
199:         if total <= 0:
200:             equal = 1.0 / len(positive_map)
201:             return {k: equal for k in positive_map}
202:         return {k: value / total for k, value in positive_map.items()}
203:
204:     @classmethod
205:     def _build_dimension_weights(
206:         cls,
207:         raw_weights: dict[str, dict[str, Any]],
208:         question_slot_lookup: dict[str, str],
209:         dimension_slot_map: dict[str, set[str]],
210:     ) -> dict[str, dict[str, float]]:
211:         dimension_weights: dict[str, dict[str, float]] = {}
212:         if raw_weights:
213:             for dim_id, weights in raw_weights.items():
214:                 resolved: dict[str, float] = {}
215:                 for qid, weight in weights.items():
216:                     slot = question_slot_lookup.get(qid, qid)
217:                     try:
218:                         resolved[slot] = float(weight)
219:                     except (TypeError, ValueError):
220:                         continue
221:                     if resolved:
222:                         dimension_weights[dim_id] = cls._normalize_weights(resolved)
223:
224:         if not dimension_weights:
```

```
225:         for dim_id, slots in dimension_slot_map.items():
226:             if not slots:
227:                 continue
228:             equal = 1.0 / len(slots)
229:             dimension_weights[dim_id] = {slot: equal for slot in slots}
230:
231:     return dimension_weights
232:
233:     @classmethod
234:     def _build_area_dimension_weights(
235:         cls,
236:         raw_weights: dict[str, dict[str, Any]],
237:         policy_areas: list[dict[str, Any]],
238:     ) -> dict[str, dict[str, float]]:
239:         area_weights: dict[str, dict[str, float]] = {}
240:         if raw_weights:
241:             for area_id, weights in raw_weights.items():
242:                 resolved: dict[str, float] = {}
243:                 for dim_id, value in weights.items():
244:                     try:
245:                         resolved[dim_id] = float(value)
246:                     except (TypeError, ValueError):
247:                         continue
248:                 if resolved:
249:                     area_weights[area_id] = cls._normalize_weights(resolved)
250:
251:         if not area_weights:
252:             for area in policy_areas:
253:                 area_id = area.get("policy_area_id") or area.get("id")
254:                 dims = area.get("dimension_ids") or []
255:                 if not area_id or not dims:
256:                     continue
257:                 equal = 1.0 / len(dims)
258:                 area_weights[area_id] = {dim: equal for dim in dims}
259:
260:     return area_weights
261:
262:     @classmethod
263:     def _build_cluster_weights(
264:         cls,
265:         raw_weights: dict[str, dict[str, Any]],
266:         clusters: list[dict[str, Any]],
267:     ) -> dict[str, dict[str, float]]:
268:         cluster_weights: dict[str, dict[str, float]] = {}
269:         if raw_weights:
270:             for cluster_id, weights in raw_weights.items():
271:                 resolved: dict[str, float] = {}
272:                 for area_id, value in weights.items():
273:                     try:
274:                         resolved[area_id] = float(value)
275:                     except (TypeError, ValueError):
276:                         continue
277:                 if resolved:
278:                     cluster_weights[cluster_id] = cls._normalize_weights(resolved)
279:
280:         if not cluster_weights:
```

```
281:         for cluster in clusters:
282:             cluster_id = cluster.get("cluster_id")
283:             area_ids = cluster.get("policy_area_ids") or []
284:             if not cluster_id or not area_ids:
285:                 continue
286:             equal = 1.0 / len(area_ids)
287:             cluster_weights[cluster_id] = {area_id: equal for area_id in area_ids}
288:
289:     return cluster_weights
290:
291:     @classmethod
292:     def _build_macro_weights(
293:         cls,
294:         raw_weights: dict[str, Any],
295:         clusters: list[dict[str, Any]],
296:     ) -> dict[str, float]:
297:         if raw_weights:
298:             resolved = {}
299:             for cluster_id, weight in raw_weights.items():
300:                 try:
301:                     resolved[cluster_id] = float(weight)
302:                 except (TypeError, ValueError):
303:                     continue
304:             normalized = cls._normalize_weights(resolved)
305:             if normalized:
306:                 return normalized
307:
308:             cluster_ids = [cluster.get("cluster_id") for cluster in clusters if cluster.get("cluster_id")]
309:             if not cluster_ids:
310:                 return {}
311:             equal = 1.0 / len(cluster_ids)
312:             return {cluster_id: equal for cluster_id in cluster_ids}
313:
314: def group_by(items: Iterable[T], key_func: Callable[[T], tuple]) -> dict[tuple, list[T]]:
315:     """
316:     Groups a sequence of items into a dictionary based on a key function.
317:
318:     This utility function iterates over a collection, applies a key function to each
319:     item, and collects items into lists, keyed by the result of the key function.
320:
321:     The key function must return a tuple. This is because dictionary keys must be
322:     hashable, and tuples are hashable whereas lists are not. Using a tuple allows
323:     for grouping by multiple attributes.
324:
325:     If the input iterable 'items' is empty, this function will return an empty
326:     dictionary.
327:
328:     Example:
329:         >>> from dataclasses import dataclass
330:         >>> @dataclass
331:         ... class Record:
332:             ...     category: str
333:             ...     value: int
334:             ...
335:         >>> data = [Record("A", 1), Record("B", 2), Record("A", 3)]
336:         >>> group_by(data, key_func=lambda r: (r.category,))
```

```
337:     {('A',): [Record(category='A', value=1), Record(category='A', value=3)],
338:      ('B',): [Record(category='B', value=2)]}
339:
340:    Args:
341:        items: An iterable of items to be grouped.
342:        key_func: A callable that accepts an item and returns a tuple to be
343:                  used as the grouping key.
344:
345:    Returns:
346:        A dictionary where keys are the result of the key function and values are
347:        lists of items belonging to that group.
348:    """
349:    grouped = defaultdict(list)
350:    for item in items:
351:        grouped[key_func(item)].append(item)
352:    return dict(grouped)
353:
354: def validate_scored_results(results: list[dict[str, Any]]) -> list[ScoredResult]:
355: """
356:     Validates a list of dictionaries and converts them to ScoredResult objects.
357:
358:     Args:
359:         results: A list of dictionaries representing scored results.
360:
361:     Returns:
362:         A list of ScoredResult objects.
363:
364:     Raises:
365:         ValidationError: If any of the dictionaries are invalid.
366:     """
367:     validated_results = []
368:     required_keys = {
369:         "question_global": int, "base_slot": str, "policy_area": str, "dimension": str,
370:         "score": float, "quality_level": str, "evidence": dict, "raw_results": dict
371:     }
372:     for i, res_dict in enumerate(results):
373:         missing_keys = set(required_keys.keys()) - set(res_dict.keys())
374:         if missing_keys:
375:             raise ValidationError(
376:                 f"Invalid ScoredResult at index {i}: missing keys {missing_keys}"
377:             )
378:         for key, expected_type in required_keys.items():
379:             if not isinstance(res_dict[key], expected_type):
380:                 raise ValidationError(
381:                     f"Invalid type for key '{key}' at index {i}. "
382:                     f"Expected {expected_type}, got {type(res_dict[key])}."
383:                 )
384:             try:
385:                 validated_results.append(ScoredResult(**res_dict))
386:             except TypeError as e:
387:                 raise ValidationError(f"Invalid ScoredResult at index {i}: {e}") from e
388:     return validated_results
389:
390: # Import canonical notation for validation
391: try:
392:     from farfan_pipeline.core.canonical_notation import get_all_dimensions, get_all_policy_areas
```

```
393:     HAS_CANONICAL_NOTATION = True
394: except ImportError:
395:     HAS_CANONICAL_NOTATION = False
396:
397: logger = logging.getLogger(__name__)
398:
399: @dataclass
400: class ScoredResult:
401:     """Represents a single, scored micro-question, forming the input for aggregation."""
402:     question_global: int
403:     base_slot: str
404:     policy_area: str
405:     dimension: str
406:     score: float
407:     quality_level: str
408:     evidence: dict[str, Any]
409:     raw_results: dict[str, Any]
410:
411: @dataclass
412: class DimensionScore:
413:     """
414:         Aggregated score for a single dimension within a policy area.
415:
416:         SOTA Extensions:
417:             - Uncertainty quantification (mean, std, CI)
418:             - Provenance tracking (DAG node ID)
419:             - Aggregation method recording
420:
421:         dimension_id: str
422:         area_id: str
423:         score: float
424:         quality_level: str
425:         contributing_questions: list[int]
426:         validation_passed: bool = True
427:         validation_details: dict[str, Any] = field(default_factory=dict)
428:
429:         # SOTA: Uncertainty quantification
430:         score_std: float = 0.0
431:         confidence_interval_95: tuple[float, float] = field(default_factory=lambda: (0.0, 0.0))
432:         epistemic_uncertainty: float = 0.0
433:         aleatoric_uncertainty: float = 0.0
434:
435:         # SOTA: Provenance tracking
436:         provenance_node_id: str = ""
437:         aggregation_method: str = "weighted_average"
438:
439: @dataclass
440: class AreaScore:
441:     """
442:         Aggregated score for a policy area, based on its constituent dimensions.
443:
444:         SOTA Extensions:
445:             - Uncertainty quantification
446:             - Provenance tracking
447:
448:         area_id: str
```

```
449:     area_name: str
450:     score: float
451:     quality_level: str
452:     dimension_scores: list[DimensionScore]
453:     validation_passed: bool = True
454:     validation_details: dict[str, Any] = field(default_factory=dict)
455:     cluster_id: str | None = None
456:
457:     # SOTA: Uncertainty quantification
458:     score_std: float = 0.0
459:     confidence_interval_95: tuple[float, float] = field(default_factory=lambda: (0.0, 0.0))
460:
461:     # SOTA: Provenance tracking
462:     provenance_node_id: str = ""
463:     aggregation_method: str = "weighted_average"
464:
465: @dataclass
466: class ClusterScore:
467:     """
468:         Aggregated score for a MESO cluster, based on its policy areas.
469:
470:         SOTA Extensions:
471:             - Uncertainty quantification
472:             - Provenance tracking
473:     """
474:     cluster_id: str
475:     cluster_name: str
476:     areas: list[str]
477:     score: float
478:     coherence: float
479:     variance: float
480:     weakest_area: str | None
481:     area_scores: list[AreaScore]
482:     validation_passed: bool = True
483:     validation_details: dict[str, Any] = field(default_factory=dict)
484:
485:     # SOTA: Uncertainty quantification
486:     score_std: float = 0.0
487:     confidence_interval_95: tuple[float, float] = field(default_factory=lambda: (0.0, 0.0))
488:
489:     # SOTA: Provenance tracking
490:     provenance_node_id: str = ""
491:     aggregation_method: str = "weighted_average"
492:
493: @dataclass
494: class MacroScore:
495:     """Represents the final, holistic macro evaluation score for the entire system."""
496:     score: float
497:     quality_level: str
498:     cross_cutting_coherence: float # Coherence across all clusters
499:     systemic_gaps: list[str]
500:     strategic_alignment: float
501:     cluster_scores: list[ClusterScore]
502:     validation_passed: bool = True
503:     validation_details: dict[str, Any] = field(default_factory=dict)
504:
```

```
505: class AggregationError(Exception):
506:     """Base exception for aggregation errors."""
507:     pass
508:
509: class ValidationError(AggregationError):
510:     """Raised when validation fails."""
511:     pass
512:
513: class WeightValidationError(ValidationError):
514:     """Raised when weight validation fails."""
515:     pass
516:
517: class ThresholdValidationError(ValidationError):
518:     """Raised when threshold validation fails."""
519:     pass
520:
521: class HermeticityValidationError(ValidationError):
522:     """Raised when hermeticity validation fails."""
523:     pass
524:
525: class CoverageError(AggregationError):
526:     """Raised when coverage requirements are not met."""
527:     pass
528:
529: class DimensionAggregator:
530:     """
531:         Aggregates micro question scores into dimension scores.
532:
533:         Responsibilities:
534:             - Aggregate 5 micro questions (Q1-Q5) per dimension
535:             - Validate weights sum to 1.0
536:             - Apply rubric thresholds
537:             - Ensure coverage (abort if insufficient)
538:             - Provide detailed logging
539:     """
540:
541:     def __init__(
542:         self,
543:         monolith: dict[str, Any] | None = None,
544:         abort_on_insufficient: bool = True,
545:         aggregation_settings: AggregationSettings | None = None,
546:         enable_sota_features: bool = True,
547:     ) -> None:
548:         """
549:             Initialize dimension aggregator.
550:
551:             Args:
552:                 monolith: Questionnaire monolith configuration (optional, required for run())
553:                 abort_on_insufficient: Whether to abort on insufficient coverage
554:                 aggregation_settings: Resolved aggregation settings
555:                 enable_sota_features: Enable SOTA features (Choquet, UQ, provenance)
556:
557:             Raises:
558:                 ValueError: If monolith is None and required for operations
559:         """
560:         self.monolith = monolith
```

```
561:         self.abort_on_insufficient = abort_on_insufficient
562:         self.aggregation_settings = aggregation_settings or AggregationSettings.from_monolith(monolith)
563:         self.dimension_group_by_keys = (
564:             self.aggregation_settings.dimension_group_by_keys or ["policy_area", "dimension"]
565:         )
566:         self.enable_sota_features = enable_sota_features
567:
568:     # Extract configuration if monolith provided
569:     if monolith is not None:
570:         self.scoring_config = monolith["blocks"]["scoring"]
571:         self.niveles = monolith["blocks"]["niveles_abstraccion"]
572:     else:
573:         self.scoring_config = None
574:         self.niveles = None
575:
576:     # SOTA: Initialize provenance DAG
577:     if self.enable_sota_features:
578:         self.provenance_dag = AggregationDAG()
579:         self.bootstrap_aggregator = BootstrapAggregator(n_samples=1000, random_seed=42)
580:         logger.info("DimensionAggregator initialized with SOTA features enabled")
581:     else:
582:         self.provenance_dag = None
583:         self.bootstrap_aggregator = None
584:         logger.info("DimensionAggregator initialized (legacy mode)")
585:
586:     # Validate canonical notation if available
587:     if HAS_CANONICAL_NOTATION:
588:         try:
589:             canonical_dims = get_all_dimensions()
590:             canonical_areas = get_all_policy_areas()
591:             logger.info(
592:                 f"Canonical notation loaded: {len(canonical_dims)} dimensions, "
593:                 f"{len(canonical_areas)} policy areas"
594:             )
595:         except Exception as e:
596:             logger.warning(f"Could not load canonical notation: {e}")
597:
598:     @calibrated_method("farfan_core.processing.aggregation.DimensionAggregator.validate_dimension_id")
599:     def validate_dimension_id(self, dimension_id: str) -> bool:
600:         """
601:             Validate dimension ID against canonical notation.
602:
603:             Args:
604:                 dimension_id: Dimension ID to validate (e.g., "DIM01")
605:
606:             Returns:
607:                 True if dimension ID is valid
608:
609:             Raises:
610:                 ValidationError: If dimension ID is invalid and abort_on_insufficient is True
611:         """
612:         if not HAS_CANONICAL_NOTATION:
613:             logger.debug("Canonical notation not available, skipping validation")
614:             return True
615:
616:         try:
```

```
617:         canonical_dims = get_all_dimensions()
618:         # Check if dimension_id is a valid code
619:         valid_codes = {info.code for info in canonical_dims.values()}
620:         if dimension_id in valid_codes:
621:             return True
622:
623:         msg = f"Invalid dimension ID: {dimension_id}. Valid codes: {sorted(valid_codes)}"
624:         logger.error(msg)
625:         if self.abort_on_insufficient:
626:             raise ValidationError(msg)
627:         return False
628:     except Exception as e:
629:         logger.warning(f"Could not validate dimension ID: {e}")
630:     return True # Don't fail if validation can't be performed
631:
632:     @calibrated_method("farfan_core.processing.aggregation.DimensionAggregator.validate_policy_area_id")
633:     def validate_policy_area_id(self, area_id: str) -> bool:
634:         """
635:             Validate policy area ID against canonical notation.
636:
637:             Args:
638:                 area_id: Policy area ID to validate (e.g., "PA01")
639:
640:             Returns:
641:                 True if policy area ID is valid
642:
643:             Raises:
644:                 ValidationError: If policy area ID is invalid and abort_on_insufficient is True
645:         """
646:         if not HAS_CANONICAL_NOTATION:
647:             logger.debug("Canonical notation not available, skipping validation")
648:             return True
649:
650:         try:
651:             canonical_areas = get_all_policy_areas()
652:             if area_id in canonical_areas:
653:                 return True
654:
655:             msg = f"Invalid policy area ID: {area_id}. Valid codes: {sorted(canonical_areas.keys())}"
656:             logger.error(msg)
657:             if self.abort_on_insufficient:
658:                 raise ValidationError(msg)
659:             return False
660:         except Exception as e:
661:             logger.warning(f"Could not validate policy area ID: {e}")
662:             return True # Don't fail if validation can't be performed
663:
664:     @calibrated_method("farfan_core.processing.aggregation.DimensionAggregator.validate_weights")
665:     def validate_weights(self, weights: list[float]) -> tuple[bool, str]:
666:         """
667:             Ensures that a list of weights sums to ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_
L582_47", 1.0) within a small tolerance.
668:
669:             Args:
670:                 weights: A list of floating-point weights.
671:
```

```
672:     Returns:
673:         A tuple containing a boolean indicating validity and a descriptive message.
674:
675:     Raises:
676:         WeightValidationError: If 'abort_on_insufficient' is True and validation fails.
677:     """
678:     if not weights:
679:         msg = "No weights provided"
680:         logger.error(msg)
681:         if self.abort_on_insufficient:
682:             raise WeightValidationError(msg)
683:         return False, msg
684:
685:     weight_sum = sum(weights)
686:     tolerance = 1e-6
687:
688:     if abs(weight_sum - ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L603_28", 1.0)) > tolerance:
689:         expected_weight = ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L604_81", 1.0)
690:         msg = f"Weight sum validation failed: sum={weight_sum:.6f}, expected={expected_weight}"
691:         logger.error(msg)
692:         if self.abort_on_insufficient:
693:             raise WeightValidationError(msg)
694:         return False, msg
695:
696:     logger.debug(f"Weight validation passed: sum={weight_sum:.6f}")
697:     return True, "Weights valid"
698:
699: def validate_coverage(
700:     self,
701:     results: list[ScoredResult],
702:     expected_count: int = 5
703: ) -> tuple[bool, str]:
704:     """
705:     Checks if the number of results meets a minimum expectation.
706:
707:     Args:
708:         results: A list of ScoredResult objects.
709:         expected_count: The minimum number of results required.
710:
711:     Returns:
712:         A tuple containing a boolean indicating validity and a descriptive message.
713:
714:     Raises:
715:         CoverageError: If 'abort_on_insufficient' is True and coverage is insufficient.
716:     """
717:     actual_count = len(results)
718:
719:     if actual_count < expected_count:
720:         msg = (
721:             f"Covariance validation failed: "
722:             f"expected {expected_count} questions, got {actual_count}"
723:         )
724:         logger.error(msg)
725:         if self.abort_on_insufficient:
726:             raise CoverageError(msg)
```

```
727:         return False, msg
728:
729:     logger.debug(f"Coverage validation passed: {actual_count}/{expected_count} questions")
730:     return True, "Coverage sufficient"
731:
732:     def calculate_weighted_average(
733:         self,
734:         scores: list[float],
735:         weights: list[float] | None = None
736:     ) -> float:
737:         """
738:             Calculates a weighted average, defaulting to an equal weighting if none provided.
739:
740:             Args:
741:                 scores: A list of scores to be averaged.
742:                 weights: An optional list of weights. If None, equal weights are assumed.
743:
744:             Returns:
745:                 The calculated weighted average.
746:
747:             Raises:
748:                 WeightValidationError: If the weights are invalid (e.g., mismatched length).
749:             """
750:
751:     if not scores:
752:         return ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L665_19", 0.0)
753:
754:     if weights is None:
755:         # Equal weights
756:         weights = [ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L669_23", 1.0) / len(scores)] * len(scores)
757:
758:     # Validate weights length matches scores length
759:     if len(weights) != len(scores):
760:         msg = (
761:             f"Weight length mismatch: {len(weights)} weights for {len(scores)} scores"
762:         )
763:         logger.error(msg)
764:         raise WeightValidationError(msg)
765:
766:     # Validate weights sum to ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L679_34", 1.0)
767:
768:     valid, msg = self.validate_weights(weights)
769:     if not valid:
770:         # If validation failed and abort_on_insufficient is False,
771:         # validate_weights already logged the error and returned False
772:         # We should raise here to avoid silent failure
773:         raise WeightValidationError(msg)
774:
775:
776:     logger.debug(
777:         f"Weighted average calculated: "
778:         f"scores={scores}, weights={weights}, result={weighted_sum:.4f}"
779:     )
780:
```

```
781:         return weighted_sum
782:
783:     def aggregate_with_sota(
784:         self,
785:         scores: list[float],
786:         weights: list[float] | None = None,
787:         method: str = "choquet",
788:         compute_uncertainty: bool = True,
789:     ) -> tuple[float, UncertaintyMetrics | None]:
790:         """
791:             SOTA aggregation with Choquet integral and uncertainty quantification.
792:
793:             This method provides:
794:             1. Non-linear aggregation via Choquet integral (captures synergies)
795:             2. Bayesian uncertainty quantification via bootstrap
796:             3. Full reproducibility with fixed random seed
797:
798:             Args:
799:                 scores: Input scores to aggregate
800:                 weights: Optional weights (default: uniform)
801:                 method: Aggregation method ("choquet" or "weighted_average")
802:                 compute_uncertainty: Whether to compute uncertainty metrics
803:
804:             Returns:
805:                 Tuple of (aggregated_score, uncertainty_metrics)
806:                 If compute_uncertainty=False, uncertainty_metrics is None
807:
808:             Raises:
809:                 ValueError: If scores is empty or method invalid
810:             """
811:             if not scores:
812:                 raise ValueError("Cannot aggregate empty score list")
813:
814:             if method == "choquet":
815:                 # Use Choquet integral for non-linear aggregation
816:                 choquet_adapter = create_default_choquet_adapter(len(scores))
817:                 score = choquet_adapter.aggregate(scores, weights)
818:                 logger.info(f"Choquet aggregation: {len(scores)} inputs \u2192 {score:.4f}")
819:             elif method == "weighted_average":
820:                 # Fall back to standard weighted average
821:                 score = self.calculate_weighted_average(scores, weights)
822:             else:
823:                 raise ValueError(f"Unknown aggregation method: {method}")
824:
825:             # Compute uncertainty if requested
826:             uncertainty = None
827:             if compute_uncertainty and self.bootstrap_aggregator:
828:                 _, uncertainty = aggregate_with_uncertainty(
829:                     scores, weights, n_bootstrap=1000, random_seed=42
830:                 )
831:                 logger.debug(
832:                     f"Uncertainty: mean={uncertainty.mean:.4f}, "
833:                     f"std={uncertainty.std:.4f}, "
834:                     f"CI95={uncertainty.confidence_interval_95}"
835:                 )
836:
```

```
837:         return score, uncertainty
838:
839:     def apply_rubric_thresholds(
840:         self,
841:         score: float,
842:         thresholds: dict[str, float] | None = None
843:     ) -> str:
844:         """
845:             Apply rubric thresholds to determine quality level.
846:
847:             Args:
848:                 score: Aggregated score (0-3 range)
849:                 thresholds: Optional threshold definitions (dict with keys: EXCELENTE, BUENO, ACEPTABLE)
850:                     Each value should be a normalized threshold (0-1 range)
851:
852:             Returns:
853:                 Quality level (EXCELENTE, BUENO, ACEPTABLE, INSUFICIENTE)
854:         """
855:         # Clamp score to valid range [0, 3]
856:         clamped_score = max(ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L714_28", 0.0), min
(3.0, score))
857:
858:         # Normalize to 0-1 range
859:         normalized_score = clamped_score / 3.0
860:
861:         # Use provided thresholds or defaults
862:         if thresholds:
863:             excellent_threshold = thresholds.get('EXCELENTE', ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights
", "auto_param_L721_62", 0.85))
864:             good_threshold = thresholds.get('BUENO', ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_
param_L722_53", 0.70))
865:             acceptable_threshold = thresholds.get('ACEPTABLE', ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weight
s", "auto_param_L723_63", 0.55))
866:             else:
867:                 excellent_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "excellent_threshold", 0.
85) # Refactored
868:                 good_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "good_threshold", 0.7) # Refac
tored
869:                 acceptable_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "acceptable_threshold",
0.55) # Refactored
870:
871:                 # Apply thresholds
872:                 if normalized_score >= excellent_threshold:
873:                     quality = "EXCELENTE"
874:                 elif normalized_score >= good_threshold:
875:                     quality = "BUENO"
876:                 elif normalized_score >= acceptable_threshold:
877:                     quality = "ACEPTABLE"
878:                 else:
879:                     quality = "INSUFICIENTE"
880:
881:                 logger.debug(
882:                     f"Rubric applied: score={score:.4f}, "
883:                     f"normalized={normalized_score:.4f}, quality={quality}"
884:                 )
885:
```

```
886:         return quality
887:
888:     def aggregate_dimension(
889:         self,
890:         scored_results: list[ScoredResult],
891:         group_by_values: dict[str, Any],
892:         weights: list[float] | None = None,
893:     ) -> DimensionScore:
894:         """
895:             Aggregate a single dimension from micro question results.
896:
897:             Args:
898:                 scored_results: List of scored results for this dimension/area.
899:                 group_by_values: Dictionary of grouping keys and their values.
900:                 weights: Optional weights for questions (defaults to equal weights).
901:
902:             Returns:
903:                 DimensionScore with aggregated score and quality level.
904:
905:             Raises:
906:                 ValidationError: If validation fails.
907:                 CoverageError: If coverage is insufficient.
908:         """
909:         dimension_id = group_by_values.get("dimension", "UNKNOWN")
910:         area_id = group_by_values.get("policy_area", "UNKNOWN")
911:         logger.info(f"Aggregating dimension {dimension_id} for area {area_id}")
912:
913:         validation_details = {}
914:
915:         # In this context, scored_results are already grouped, so we can use them directly.
916:         dim_results = scored_results
917:
918:         expected_count = self._expected_question_count(area_id, dimension_id)
919:
920:         # Validate coverage
921:         try:
922:             coverage_valid, coverage_msg = self.validate_coverage(
923:                 dim_results,
924:                 expected_count=expected_count or 5,
925:             )
926:             validation_details["coverage"] = {
927:                 "valid": coverage_valid,
928:                 "message": coverage_msg,
929:                 "count": len(dim_results)
930:             }
931:         except CoverageError as e:
932:             logger.error(f"Cov  
erage validation failed for {dimension_id}/{area_id}: {e}")
933:             # Return minimal score if aborted
934:             return DimensionScore(
935:                 dimension_id=dimension_id,
936:                 area_id=area_id,
937:                 score=ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L795_22", 0.0),
938:                 quality_level="INSUFICIENTE",
939:                 contributing_questions=[],
940:                 validation_passed=False,
941:                 validation_details={"error": str(e), "type": "coverage"}
```

```
942:         )
943:
944:     if not dim_results:
945:         logger.warning(f"No results for dimension {dimension_id}/{area_id}")
946:         return DimensionScore(
947:             dimension_id=dimension_id,
948:             area_id=area_id,
949:             score=ParameterLoaderV2.get("farfan_core.processing.aggregation.DimensionAggregator.validate_weights", "auto_param_L807_22", 0.0),
950:             quality_level="INSUFICIENTE",
951:             contributing_questions=[],
952:             validation_passed=False,
953:             validation_details={"error": "No results", "type": "empty"}
954:         )
955:
956:     # Extract scores
957:     scores = [r.score for r in dim_results]
958:     question_ids = [r.question_global for r in dim_results]
959:
960:     # Calculate weighted average with SOTA features
961:     resolved_weights = weights or self._resolve_dimension_weights(dimension_id, dim_results)
962:
963:     # SOTA: Use Choquet + uncertainty if enabled
964:     if self.enable_sota_features and len(scores) >= 3:
965:         try:
966:             avg_score, uncertainty = self.aggregate_with_sota(
967:                 scores,
968:                 resolved_weights,
969:                 method="choquet",
970:                 compute_uncertainty=True,
971:             )
972:             validation_details["aggregation"] = {
973:                 "method": "choquet",
974:                 "uncertainty": uncertainty.to_dict() if uncertainty else None,
975:             }
976:         except Exception as e:
977:             logger.warning(f"SOTA aggregation failed, falling back to standard: {e}")
978:             avg_score = self.calculate_weighted_average(scores, resolved_weights)
979:             uncertainty = None
980:             validation_details["aggregation"] = {"method": "weighted_average", "fallback": True}
981:     else:
982:         # Standard aggregation
983:         avg_score = self.calculate_weighted_average(scores, resolved_weights)
984:         uncertainty = None
985:         validation_details["aggregation"] = {"method": "weighted_average"}
986:
987:     validation_details["weights"] = {
988:         "valid": True,
989:         "weights": resolved_weights if resolved_weights else "equal",
990:         "score": avg_score
991:     }
992:
993:     # Apply rubric thresholds
994:     quality_level = self.apply_rubric_thresholds(avg_score)
995:     validation_details["rubric"] = {
996:         "score": avg_score,
997:         "quality_level": quality_level
998:     }
```

```

998:         }
999:     validation_details["score_max"] = 3.0
1000:
1001:     # SOTA: Add provenance tracking
1002:     provenance_node_id = f"DIM_{dimension_id}_{area_id}"
1003:     if self.enable_sota_features and self.provenance_dag:
1004:         # Add dimension node
1005:         dim_node = ProvenanceNode(
1006:             node_id=provenance_node_id,
1007:             level="dimension",
1008:             score=avg_score,
1009:             quality_level=quality_level,
1010:             metadata={
1011:                 "dimension_id": dimension_id,
1012:                 "area_id": area_id,
1013:                 "n_questions": len(question_ids),
1014:             },
1015:         )
1016:         self.provenance_dag.add_node(dim_node)
1017:
1018:     # Add aggregation edges from questions to dimension
1019:     question_node_ids = [f"Q{qid:03d}" for qid in question_ids]
1020:     for qid_str, qid in zip(question_node_ids, question_ids):
1021:         # Add question node if not exists
1022:         if qid_str not in self.provenance_dag.nodes:
1023:             q_node = ProvenanceNode(
1024:                 node_id=qid_str,
1025:                 level="micro",
1026:                 score=scores[question_ids.index(qid)],
1027:                 quality_level="UNKNOWN",
1028:             )
1029:             self.provenance_dag.add_node(q_node)
1030:
1031:     # Record aggregation operation
1032:     self.provenance_dag.add_aggregation_edge(
1033:         source_ids=question_node_ids,
1034:         target_id=provenance_node_id,
1035:         operation="choquet" if self.enable_sota_features else "weighted_average",
1036:         weights=resolved_weights or [1.0 / len(question_ids)] * len(question_ids),
1037:         metadata={"dimension": dimension_id, "area": area_id},
1038:     )
1039:
1040:     logger.debug(f"Provenance recorded: {len(question_node_ids)} questions \u2192 {provenance_node_id}")
1041:
1042:     logger.info(
1043:         f"\u2192 Dimension {dimension_id}/{area_id}: "
1044:         f"score={avg_score:.4f}, quality={quality_level}"
1045:         + (f", std={uncertainty.std:.4f}" if uncertainty else ""))
1046: )
1047:
1048:     return DimensionScore(
1049:         dimension_id=dimension_id,
1050:         area_id=area_id,
1051:         score=avg_score,
1052:         quality_level=quality_level,
1053:         contributing_questions=question_ids,

```

```
1054:         validation_passed=True,
1055:         validation_details=validation_details,
1056:         # SOTA fields
1057:         score_std=uncertainty.std if uncertainty else 0.0,
1058:         confidence_interval_95=uncertainty.confidence_interval_95 if uncertainty else (0.0, 0.0),
1059:         epistemic_uncertainty=uncertainty.epistemic_uncertainty if uncertainty else 0.0,
1060:         aleatoric_uncertainty=uncertainty.aleatoric_uncertainty if uncertainty else 0.0,
1061:         provenance_node_id=provenance_node_id if self.enable_sota_features else "",
1062:         aggregation_method="choquet" if (self.enable_sota_features and len(scores) >= 3) else "weighted_average",
1063:     )
1064:
1065:     def run(
1066:         self,
1067:         scored_results: list[ScoredResult],
1068:         group_by_keys: list[str]
1069:     ) -> list[DimensionScore]:
1070:         """
1071:             Run the dimension aggregation process.
1072:
1073:             Args:
1074:                 scored_results: List of all scored results.
1075:                 group_by_keys: List of keys to group by.
1076:
1077:             Returns:
1078:                 A list of DimensionScore objects.
1079:         """
1080:         def key_func(r):
1081:             return tuple(getattr(r, key) for key in group_by_keys)
1082:         grouped_results = group_by(scored_results, key_func)
1083:
1084:         dimension_scores = []
1085:         for group_key, results in grouped_results.items():
1086:             group_by_values = dict(zip(group_by_keys, group_key, strict=False))
1087:             score = self.aggregate_dimension(results, group_by_values)
1088:             dimension_scores.append(score)
1089:
1090:         return dimension_scores
1091:
1092:     @calibrated_method("farfan_core.processing.aggregation.DimensionAggregator._expected_question_count")
1093:     def _expected_question_count(self, area_id: str, dimension_id: str) -> int | None:
1094:         if not self.aggregation_settings.dimension_expected_counts:
1095:             return None
1096:         return self.aggregation_settings.dimension_expected_counts.get((area_id, dimension_id))
1097:
1098:     def _resolve_dimension_weights(
1099:         self,
1100:         dimension_id: str,
1101:         dim_results: list[ScoredResult],
1102:     ) -> list[float] | None:
1103:         mapping = self.aggregation_settings.dimension_question_weights.get(dimension_id)
1104:         if not mapping:
1105:             return None
1106:
1107:         weights: list[float] = []
1108:         for result in dim_results:
1109:             slot = result.base_slot
```

```
1110:         weight = mapping.get(slot)
1111:         if weight is None:
1112:             logger.debug(
1113:                 "Missing weight for slot %s in dimension %s \u200\223 falling back to equal weights",
1114:                 slot,
1115:                 dimension_id,
1116:             )
1117:         return None
1118:     weights.append(weight)
1119:
1120:     total = sum(weights)
1121:     if total <= 0:
1122:         return None
1123:     return [w / total for w in weights]
1124:
1125: def run_aggregation_pipeline(
1126:     scored_results: list[dict[str, Any]],
1127:     monolith: dict[str, Any],
1128:     abort_on_insufficient: bool = True
1129: ) -> list[ClusterScore]:
1130: """
1131: Orchestrates the end-to-end aggregation pipeline.
1132:
1133: This function provides a high-level entry point to the aggregation system,
1134: demonstrating the sequential wiring of the aggregator components. It ensures
1135: that data flows from raw scored results through dimension, area, and
1136: finally cluster aggregation in a controlled and validated manner.
1137:
1138: Note on Parallelization: This implementation is sequential. For very large
1139: datasets, the 'group_by' operations in each aggregator's 'run' method
1140: could be parallelized (e.g., using 'concurrent.futures') to process
1141: independent groups concurrently.
1142:
1143: Args:
1144:     scored_results: A list of dictionaries, each representing a raw scored result.
1145:     monolith: The central monolith configuration object.
1146:     abort_on_insufficient: If True, the pipeline will stop on validation errors.
1147:
1148: Returns:
1149:     A list of aggregated ClusterScore objects.
1150: """
1151: # 1. Input Validation (Pre-flight check)
1152: validated_scored_results = validate_scored_results(scored_results)
1153:
1154: aggregation_settings = AggregationSettings.from_monolith(monolith)
1155:
1156: # 2. FASE 4: Dimension Aggregation
1157: dim_aggregator = DimensionAggregator(
1158:     monolith,
1159:     abort_on_insufficient,
1160:     aggregation_settings=aggregation_settings,
1161: )
1162: dimension_scores = dim_aggregator.run(
1163:     validated_scored_results,
1164:     group_by_keys=dim_aggregator.dimension_group_by_keys,
1165: )
```

```
1166:  
1167:     # 3. FASE 5: Area Policy Aggregation  
1168:     area_aggregator = AreaPolicyAggregator(  
1169:         monolith,  
1170:         abort_on_insufficient,  
1171:         aggregation_settings=aggregation_settings,  
1172:     )  
1173:     area_scores = area_aggregator.run(  
1174:         dimension_scores,  
1175:         group_by_keys=area_aggregator.area_group_by_keys,  
1176:     )  
1177:  
1178:     # 4. FASE 6: Cluster Aggregation  
1179:     cluster_aggregator = ClusterAggregator(  
1180:         monolith,  
1181:         abort_on_insufficient,  
1182:         aggregation_settings=aggregation_settings,  
1183:     )  
1184:     cluster_definitions = monolith["blocks"]["niveles_abstraccion"]["clusters"]  
1185:     cluster_scores = cluster_aggregator.run(  
1186:         area_scores,  
1187:         cluster_definitions  
1188:     )  
1189:  
1190:     return cluster_scores  
1191:  
1192: def run(  
1193:     self,  
1194:     scored_results: list[ScoredResult],  
1195:     group_by_keys: list[str]  
1196: ) -> list[DimensionScore]:  
1197:     """  
1198:         Run the dimension aggregation process.  
1199:  
1200:     Args:  
1201:         scored_results: List of all scored results.  
1202:         group_by_keys: List of keys to group by.  
1203:  
1204:     Returns:  
1205:         A list of DimensionScore objects.  
1206:     """  
1207:     def key_func(r):  
1208:         return tuple(getattr(r, key) for key in group_by_keys)  
1209:     grouped_results = group_by(scored_results, key_func)  
1210:  
1211:     dimension_scores = []  
1212:     for group_key, results in grouped_results.items():  
1213:         group_by_values = dict(zip(group_by_keys, group_key, strict=False))  
1214:         score = self.aggregate_dimension(results, group_by_values)  
1215:         dimension_scores.append(score)  
1216:  
1217:     return dimension_scores  
1218:  
1219: class AreaPolicyAggregator:  
1220:     """  
1221:         Aggregates dimension scores into policy area scores.
```

```
1222:
1223:     Responsibilities:
1224:         - Aggregate 6 dimension scores per policy area
1225:         - Validate dimension completeness
1226:         - Apply area-level rubric thresholds
1227:         - Ensure hermeticity (no dimension overlap)
1228:     """
1229:
1230:     def __init__(
1231:         self,
1232:             monolith: dict[str, Any] | None = None,
1233:             abort_on_insufficient: bool = True,
1234:             aggregation_settings: AggregationSettings | None = None,
1235:         ) -> None:
1236:         """
1237:             Initialize area aggregator.
1238:
1239:         Args:
1240:             monolith: Questionnaire monolith configuration (optional, required for run())
1241:             abort_on_insufficient: Whether to abort on insufficient coverage
1242:
1243:         Raises:
1244:             ValueError: If monolith is None and required for operations
1245:         """
1246:         self.monolith = monolith
1247:         self.abort_on_insufficient = abort_on_insufficient
1248:         self.aggregation_settings = aggregation_settings or AggregationSettings.from_monolith(monolith)
1249:         self.area_group_by_keys = self.aggregation_settings.area_group_by_keys or ["area_id"]
1250:
1251:         # Extract configuration if monolith provided
1252:         if monolith is not None:
1253:             self.scoring_config = monolith["blocks"]["scoring"]
1254:             self.niveles = monolith["blocks"]["niveles_abstraccion"]
1255:             self.policy_areas = self.niveles["policy_areas"]
1256:             self.dimensions = self.niveles["dimensions"]
1257:         else:
1258:             self.scoring_config = None
1259:             self.niveles = None
1260:             self.policy_areas = None
1261:             self.dimensions = None
1262:
1263:             logger.info("AreaPolicyAggregator initialized")
1264:
1265:     def validate_hermeticity(
1266:         self,
1267:             dimension_scores: list[DimensionScore],
1268:             area_id: str
1269:         ) -> tuple[bool, str]:
1270:         """
1271:             Validate hermeticity (no dimension overlap/gaps).
1272:             Uses scoped validation based on policy_area.dimension_ids from monolith.
1273:
1274:         Args:
1275:             dimension_scores: List of dimension scores for the area
1276:             area_id: Policy area ID
1277:
```

```
1278:     Returns:
1279:         Tuple of (is_valid, message)
1280:
1281:     Raises:
1282:         HermeticityValidationError: If hermeticity is violated
1283:         """
1284:     # Get expected dimensions for this specific policy area
1285:     area_def = next(
1286:         (a for a in self.policy_areas if a["policy_area_id"] == area_id),
1287:         None
1288:     )
1289:
1290:     if area_def and "dimension_ids" in area_def:
1291:         expected_dimension_ids = set(area_def["dimension_ids"])
1292:     else:
1293:         # Fallback to all global dimensions if not specified
1294:         expected_dimension_ids = {d["dimension_id"] for d in self.dimensions}
1295:
1296:     actual_dimension_ids = {d.dimension_id for d in dimension_scores}
1297:     len(expected_dimension_ids)
1298:     len(dimension_scores)
1299:
1300:     # Check for missing dimensions
1301:     missing_dims = expected_dimension_ids - actual_dimension_ids
1302:     if missing_dims:
1303:         msg = (
1304:             f"Hermeticity violation for area {area_id}: "
1305:             f"missing dimensions {missing_dims}"
1306:         )
1307:         logger.error(msg)
1308:         if self.abort_on_insufficient:
1309:             raise HermeticityValidationError(msg)
1310:         return False, msg
1311:
1312:     # Check for unexpected dimensions
1313:     extra_dims = actual_dimension_ids - expected_dimension_ids
1314:     if extra_dims:
1315:         msg = (
1316:             f"Hermeticity violation for area {area_id}: "
1317:             f"unexpected dimensions {extra_dims}"
1318:         )
1319:         logger.error(msg)
1320:         if self.abort_on_insufficient:
1321:             raise HermeticityValidationError(msg)
1322:         return False, msg
1323:
1324:     # Check for duplicate dimensions
1325:     dimension_ids = [d.dimension_id for d in dimension_scores]
1326:     if len(dimension_ids) != len(set(dimension_ids)):
1327:         msg = f"Hermeticity violation for area {area_id}: duplicate dimensions found"
1328:         logger.error(msg)
1329:         if self.abort_on_insufficient:
1330:             raise HermeticityValidationError(msg)
1331:         return False, msg
1332:
1333:     logger.debug(f"Hermeticity validation passed for area {area_id}")
```

```
1334:         return True, "Hermeticity validated"
1335:
1336:     @calibrated_method("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores")
1337:     def normalize_scores(self, dimension_scores: list[DimensionScore]) -> list[float]:
1338:         """
1339:             Normalize dimension scores to 0-1 range.
1340:
1341:             Args:
1342:                 dimension_scores: List of dimension scores
1343:
1344:             Returns:
1345:                 List of normalized scores
1346:             """
1347:             normalized = []
1348:             for d in dimension_scores:
1349:                 # Extract max_expected from validation_details or default to 3.0
1350:                 max_expected = d.validation_details.get('score_max', 3.0) if d.validation_details else 3.0
1351:                 normalized.append(max(ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1148_34", 0
1352:                 .0), min(max_expected, d.score)) / max_expected)
1353:             logger.debug(f"Scores normalized: {normalized}")
1354:             return normalized
1355:
1356:     def apply_rubric_thresholds(
1357:         self,
1358:         score: float,
1359:         thresholds: dict[str, float] | None = None
1360:     ) -> str:
1361:         """
1362:             Apply area-level rubric thresholds.
1363:
1364:             Args:
1365:                 score: Aggregated score (0-3 range)
1366:                 thresholds: Optional threshold definitions (dict with keys: EXCELENTE, BUENO, ACEPTABLE)
1367:                     Each value should be a normalized threshold (0-1 range)
1368:
1369:             Returns:
1370:                 Quality level (EXCELENTE, BUENO, ACEPTABLE, INSUFICIENTE)
1371:             """
1372:             # Clamp score to valid range [0, 3]
1373:             clamped_score = max(ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1170_28", 0.0), m
in(3.0, score))
1374:
1375:             # Normalize to 0-1 range
1376:             normalized_score = clamped_score / 3.0
1377:
1378:             # Use provided thresholds or defaults
1379:             if thresholds:
1380:                 excellent_threshold = thresholds.get('EXCELENTE', ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores",
1381:                     "auto_param_L1177_62", 0.85))
1382:                 good_threshold = thresholds.get('BUENO', ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto
1383:                     _param_L1178_53", 0.70))
1384:                 acceptable_threshold = thresholds.get('ACEPTABLE', ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores",
1385:                     "auto_param_L1179_63", 0.55))
1386:             else:
1387:                 excellent_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "excellent_threshold", 0
```

```

.85) # Refactored
1385:     good_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "good_threshold", 0.7) # Refa
ctored
1386:     acceptable_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "acceptable_threshold",
0.55) # Refactored
1387:
1388:     # Apply thresholds
1389:     if normalized_score >= excellent_threshold:
1390:         quality = "EXCELENTE"
1391:     elif normalized_score >= good_threshold:
1392:         quality = "BUENO"
1393:     elif normalized_score >= acceptable_threshold:
1394:         quality = "ACEPTABLE"
1395:     else:
1396:         quality = "INSUFICIENTE"
1397:
1398:     logger.debug(
1399:         f"Area rubric applied: score={score:.4f}, "
1400:         f"normalized={normalized_score:.4f}, quality={quality}"
1401:     )
1402:
1403:     return quality
1404:
1405: def aggregate_area(
1406:     self,
1407:     dimension_scores: list[DimensionScore],
1408:     group_by_values: dict[str, Any],
1409:     weights: list[float] | None = None,
1410: ) -> AreaScore:
1411:     """
1412:     Aggregate a single policy area from dimension scores.
1413:
1414:     Args:
1415:         dimension_scores: List of dimension scores for this area.
1416:         group_by_values: Dictionary of grouping keys and their values.
1417:         weights: Optional list of weights for dimension scores.
1418:
1419:     Returns:
1420:         AreaScore with aggregated score and quality level.
1421:
1422:     Raises:
1423:         ValidationError: If validation fails.
1424:     """
1425:     area_id = group_by_values.get("area_id", "UNKNOWN")
1426:     logger.info(f"Aggregating policy area {area_id}")
1427:
1428:     validation_details = {}
1429:
1430:     # The dimension_scores are already grouped.
1431:     area_dim_scores = dimension_scores
1432:
1433:     # Validate hermeticity
1434:     try:
1435:         hermetic_valid, hermetic_msg = self.validate_hermeticity(area_dim_scores, area_id)
1436:         validation_details["hermeticity"] = {
1437:             "valid": hermetic_valid,

```

```

1438:             "message": hermetic_msg,
1439:             "dimension_count": len(area_dim_scores)
1440:         }
1441:     except HermeticityValidationError as e:
1442:         logger.error(f"Hermeticity validation failed for area {area_id}: {e}")
1443:         # Get area name
1444:         area_name = next(
1445:             (a["i18n"]["keys"]["label_es"] for a in self.policy_areas
1446:              if a["policy_area_id"] == area_id),
1447:             area_id
1448:         )
1449:     return AreaScore(
1450:         area_id=area_id,
1451:         area_name=area_name,
1452:         score=ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1249_22", 0.0),
1453:         quality_level="INSUFICIENTE",
1454:         dimension_scores=[],
1455:         validation_passed=False,
1456:         validation_details={"error": str(e), "type": "hermeticity"}
1457:     )
1458:
1459:     if not area_dim_scores:
1460:         logger.warning(f"No dimension scores for area {area_id}")
1461:         area_name = next(
1462:             (a["i18n"]["keys"]["label_es"] for a in self.policy_areas
1463:              if a["policy_area_id"] == area_id),
1464:             area_id
1465:         )
1466:     return AreaScore(
1467:         area_id=area_id,
1468:         area_name=area_name,
1469:         score=ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1266_22", 0.0),
1470:         quality_level="INSUFICIENTE",
1471:         dimension_scores=[],
1472:         validation_passed=False,
1473:         validation_details={"error": "No dimensions", "type": "empty"}
1474:     )
1475:
1476:     # Normalize scores
1477:     normalized = self.normalize_scores(area_dim_scores)
1478:     validation_details["normalization"] = {
1479:         "original": [d.score for d in area_dim_scores],
1480:         "normalized": normalized
1481:     }
1482:
1483:     # Calculate weighted average score
1484:     scores = [d.score for d in area_dim_scores]
1485:     resolved_weights = weights or self._resolve_area_weights(area_id, area_dim_scores)
1486:     avg_score = DimensionAggregator().calculate_weighted_average(scores, weights=resolved_weights)
1487:
1488:     # Apply rubric thresholds
1489:     quality_level = self.apply_rubric_thresholds(avg_score)
1490:     validation_details["rubric"] = {
1491:         "score": avg_score,
1492:         "quality_level": quality_level
1493:     }

```

```
1494:  
1495:     # Get area name  
1496:     area_name = next(  
1497:         (a["i18n"]["keys"]["label_es"] for a in self.policy_areas  
1498:             if a["policy_area_id"] == area_id),  
1499:         area_id  
1500:     )  
1501:  
1502:     logger.info(  
1503:         f"\u234\u223 Policy area {area_id} ({area_name}): "  
1504:         f"score={avg_score:.4f}, quality={quality_level}"  
1505:     )  
1506:  
1507:     return AreaScore(  
1508:         area_id=area_id,  
1509:         area_name=area_name,  
1510:         score=avg_score,  
1511:         quality_level=quality_level,  
1512:         dimension_scores=area_dim_scores,  
1513:         validation_passed=True,  
1514:         validation_details=validation_details  
1515:     )  
1516:  
1517:     def run(  
1518:         self,  
1519:         dimension_scores: list[DimensionScore],  
1520:         group_by_keys: list[str]  
1521:     ) -> list[AreaScore]:  
1522:         """  
1523:             Run the area aggregation process.  
1524:  
1525:             Args:  
1526:                 dimension_scores: List of all dimension scores.  
1527:                 group_by_keys: List of keys to group by.  
1528:  
1529:             Returns:  
1530:                 A list of AreaScore objects.  
1531:         """  
1532:         def key_func(d):  
1533:             return tuple(getattr(d, key) for key in group_by_keys)  
1534:         grouped_scores = group_by(dimension_scores, key_func)  
1535:  
1536:         area_scores = []  
1537:         for group_key, scores in grouped_scores.items():  
1538:             group_by_values = dict(zip(group_by_keys, group_key, strict=False))  
1539:             score = self.aggregate_area(scores, group_by_values, weights=None)  
1540:             area_scores.append(score)  
1541:  
1542:         return area_scores  
1543:  
1544:     def _resolve_area_weights(  
1545:         self,  
1546:         area_id: str,  
1547:         dimension_scores: list[DimensionScore],  
1548:     ) -> list[float] | None:  
1549:         mapping = self.aggregation_settings.policy_area_dimension_weights.get(area_id)
```

```

1550:         if not mapping:
1551:             return None
1552:
1553:         weights: list[float] = []
1554:         for dim_score in dimension_scores:
1555:             weight = mapping.get(dim_score.dimension_id)
1556:             if weight is None:
1557:                 logger.debug(
1558:                     "Missing weight for dimension %s in area %s \u200\223 falling back to equal weights",
1559:                     dim_score.dimension_id,
1560:                     area_id,
1561:                 )
1562:             return None
1563:         weights.append(weight)
1564:
1565:     total = sum(weights)
1566:     if total <= 0:
1567:         return None
1568:     return [w / total for w in weights]
1569:
1570: class ClusterAggregator:
1571:     """
1572:     Aggregates policy area scores into cluster scores (MESO level).
1573:
1574:     Responsibilities:
1575:     - Aggregate multiple area scores per cluster
1576:     - Apply cluster-specific weights
1577:     - Calculate coherence metrics
1578:     - Validate cluster hermeticity
1579:     """
1580:
1581: PENALTY_WEIGHT = ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "PENALTY_WEIGHT", 0.3)
1582: MAX_SCORE = 3.0
1583:
1584: def __init__(
1585:     self,
1586:     monolith: dict[str, Any] | None = None,
1587:     abort_on_insufficient: bool = True,
1588:     aggregation_settings: AggregationSettings | None = None,
1589: ) -> None:
1590:     """
1591:     Initialize cluster aggregator.
1592:
1593:     Args:
1594:         monolith: Questionnaire monolith configuration (optional, required for run())
1595:         abort_on_insufficient: Whether to abort on insufficient coverage
1596:
1597:     Raises:
1598:         ValueError: If monolith is None and required for operations
1599:     """
1600:     self.monolith = monolith
1601:     self.abort_on_insufficient = abort_on_insufficient
1602:     self.aggregation_settings = aggregation_settings or AggregationSettings.from_monolith(monolith)
1603:     self.cluster_group_by_keys = self.aggregation_settings.cluster_group_by_keys or ["cluster_id"]
1604:
1605:     # Extract configuration if monolith provided

```

```
1606:         if monolith is not None:
1607:             self.scoring_config = monolith["blocks"]["scoring"]
1608:             self.niveles = monolith["blocks"]["niveles_abstraccion"]
1609:             self.clusters = self.niveles["clusters"]
1610:         else:
1611:             self.scoring_config = None
1612:             self.niveles = None
1613:             self.clusters = None
1614:
1615:         logger.info("ClusterAggregator initialized")
1616:
1617:     def validate_cluster_hermeticity(
1618:         self,
1619:         cluster_def: dict[str, Any],
1620:         area_scores: list[AreaScore]
1621:     ) -> tuple[bool, str]:
1622:         """
1623:             Validate cluster hermeticity.
1624:
1625:             Args:
1626:                 cluster_def: Cluster definition from monolith
1627:                 area_scores: List of area scores for this cluster
1628:
1629:             Returns:
1630:                 Tuple of (is_valid, message)
1631:
1632:             Raises:
1633:                 HermeticityValidationError: If hermeticity is violated
1634:         """
1635:         expected_areas = cluster_def.get("policy_area_ids", [])
1636:         actual_areas = [a.area_id for a in area_scores]
1637:
1638:         # Check for duplicate areas
1639:         if len(actual_areas) != len(set(actual_areas)):
1640:             msg = (
1641:                 f"Cluster hermeticity violation: "
1642:                 f"duplicate areas found for cluster {cluster_def['cluster_id']}"
1643:             )
1644:             logger.error(msg)
1645:             if self.abort_on_insufficient:
1646:                 raise HermeticityValidationError(msg)
1647:             return False, msg
1648:
1649:         # Check that all expected areas are present
1650:         missing_areas = set(expected_areas) - set(actual_areas)
1651:         if missing_areas:
1652:             msg = (
1653:                 f"Cluster hermeticity violation: "
1654:                 f"missing areas {missing_areas} for cluster {cluster_def['cluster_id']}"
1655:             )
1656:             logger.error(msg)
1657:             if self.abort_on_insufficient:
1658:                 raise HermeticityValidationError(msg)
1659:             return False, msg
1660:
1661:         # Check for unexpected areas
```

```
1662:     extra_areas = set(actual_areas) - set(expected_areas)
1663:     if extra_areas:
1664:         msg = (
1665:             f"Cluster hermeticity violation: "
1666:             f"unexpected areas {extra_areas} for cluster {cluster_def['cluster_id']}"
1667:         )
1668:         logger.error(msg)
1669:         if self.abort_on_insufficient:
1670:             raise HermeticityValidationError(msg)
1671:         return False, msg
1672:
1673:     logger.debug(f"Cluster hermeticity validated for {cluster_def['cluster_id']}")
1674:     return True, "Cluster hermeticity validated"
1675:
1676: def apply_cluster_weights(
1677:     self,
1678:     area_scores: list[AreaScore],
1679:     weights: list[float] | None = None
1680: ) -> float:
1681: """
1682:     Apply cluster-specific weights to area scores.
1683:
1684:     Args:
1685:         area_scores: List of area scores
1686:         weights: Optional weights (defaults to equal weights)
1687:
1688:     Returns:
1689:         Weighted average score
1690:
1691:     Raises:
1692:         WeightValidationError: If weights validation fails
1693: """
1694: scores = [a.score for a in area_scores]
1695:
1696: if weights is None:
1697:     # Equal weights
1698:     weights = [ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1495_23", 1.0) / len(scores)] * len(scores)
1699:
1700:     # Validate weights length matches scores length
1701:     if len(weights) != len(scores):
1702:         msg = (
1703:             f"Cluster weight length mismatch: "
1704:             f"{len(weights)} weights for {len(scores)} area scores"
1705:         )
1706:         logger.error(msg)
1707:         if self.abort_on_insufficient:
1708:             raise WeightValidationError(msg)
1709:
1710:     # Validate weights sum to ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1507_34", 1
1711:     weight_sum = sum(weights)
1712:     tolerance = 1e-6
1713:     if abs(weight_sum - ParameterLoaderV2.get("farfan_core.processing.aggregation.AreaPolicyAggregator.normalize_scores", "auto_param_L1510_28", 1.0)) >
tolerance:
1714:         msg = f"Cluster weight validation failed: sum={weight_sum:.6f}"
```

```
1715:         logger.error(msg)
1716:         if self.abort_on_insufficient:
1717:             raise WeightValidationError(msg)
1718:
1719:         # Calculate weighted average
1720:         weighted_avg = sum(s * w for s, w in zip(scores, weights, strict=False))
1721:
1722:         logger.debug(
1723:             f"Cluster weights applied: scores={scores}, "
1724:             f"weights={weights}, result={weighted_avg:.4f}"
1725:         )
1726:
1727:         return weighted_avg
1728:
1729:     @calibrated_method("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence")
1730:     def analyze_coherence(self, area_scores: list[AreaScore]) -> float:
1731:         """
1732:             Analyze cluster coherence.
1733:
1734:             Coherence is measured as the inverse of standard deviation.
1735:             Higher coherence means scores are more consistent.
1736:
1737:             Args:
1738:                 area_scores: List of area scores
1739:
1740:             Returns:
1741:                 Coherence value (0-1, where 1 is perfect coherence)
1742: """
1743:     scores = [a.score for a in area_scores]
1744:
1745:     if len(scores) <= 1:
1746:         return ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1543_19", 1.0)
1747:
1748:     # Calculate mean
1749:     mean = sum(scores) / len(scores)
1750:
1751:     # Calculate standard deviation
1752:     variance = sum((s - mean) ** 2 for s in scores) / len(scores)
1753:     std_dev = variance ** ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1550_30", 0.5)
1754:
1755:     # Convert to coherence (inverse relationship)
1756:     # Normalize by max possible std dev (3.0 for 0-3 range)
1757:     max_std = 3.0
1758:     coherence = max(ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1555_24", 0.0), ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1555_29", 1.0) - (std_dev / max_std))
1759:
1760:     logger.debug(
1761:         f"Coherence analysis: mean={mean:.4f}, "
1762:             f"std_dev={std_dev:.4f}, coherence={coherence:.4f}"
1763:     )
1764:
1765:     return coherence
1766:
1767:     def aggregate_cluster(
1768:         self,
1769:         area_scores: list[AreaScore],
```

```
1770:         group_by_values: dict[str, Any],
1771:         weights: list[float] | None = None,
1772:     ) -> ClusterScore:
1773:         """
1774:             Aggregate a single MESO cluster from area scores.
1775:
1776:             Args:
1777:                 area_scores: List of area scores for this cluster.
1778:                 group_by_values: Dictionary of grouping keys and their values.
1779:                 weights: Optional cluster-specific weights.
1780:
1781:             Returns:
1782:                 ClusterScore with aggregated score and coherence.
1783:
1784:             Raises:
1785:                 ValidationError: If validation fails.
1786:             """
1787:             cluster_id = group_by_values.get("cluster_id", "UNKNOWN")
1788:             logger.info(f"Aggregating cluster {cluster_id}")
1789:
1790:             validation_details = {}
1791:
1792:             # Get cluster definition
1793:             cluster_def = next(
1794:                 (c for c in self.clusters if c["cluster_id"] == cluster_id), None
1795:             )
1796:
1797:             if not cluster_def:
1798:                 logger.error(f"Cluster definition not found: {cluster_id}")
1799:                 return ClusterScore(
1800:                     cluster_id=cluster_id,
1801:                     cluster_name=cluster_id,
1802:                     areas=[],
1803:                     score=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1600_22", 0.0),
1804:                     coherence=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1601_26", 0.0),
1805:                     variance=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1602_25", 0.0),
1806:                     weakest_area=None,
1807:                     area_scores=[],
1808:                     validation_passed=False,
1809:                     validation_details={"error": "Definition not found", "type": "config"},
1810:                 )
1811:
1812:             cluster_name = cluster_def["i18n"]["keys"]["label_es"]
1813:             expected_areas = cluster_def["policy_area_ids"]
1814:
1815:             # The area_scores are already grouped.
1816:             cluster_area_scores = area_scores
1817:
1818:             # Validate hermeticity
1819:             try:
1820:                 hermetic_valid, hermetic_msg = self.validate_cluster_hermeticity(
1821:                     cluster_def,
1822:                     cluster_area_scores
1823:                 )
1824:                 validation_details["hermeticity"] = {
1825:                     "valid": hermetic_valid,
```

```
1826:         "message": hermetic_msg
1827:     }
1828: except HermeticityValidationError as e:
1829:     logger.error(f"Cluster hermeticity validation failed: {e}")
1830:     return ClusterScore(
1831:         cluster_id=cluster_id,
1832:         cluster_name=cluster_name,
1833:         areas=expected_areas,
1834:         score=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1631_22", 0.0),
1835:         coherence=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1632_26", 0.0),
1836:         variance=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1633_25", 0.0),
1837:         weakest_area=None,
1838:         area_scores=[],
1839:         validation_passed=False,
1840:         validation_details={"error": str(e), "type": "hermeticity"}
1841:     )
1842:
1843: if not cluster_area_scores:
1844:     logger.warning(f"No area scores for cluster {cluster_id}")
1845:     return ClusterScore(
1846:         cluster_id=cluster_id,
1847:         cluster_name=cluster_name,
1848:         areas=expected_areas,
1849:         score=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1646_22", 0.0),
1850:         coherence=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1647_26", 0.0),
1851:         variance=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1648_25", 0.0),
1852:         weakest_area=None,
1853:         area_scores=[],
1854:         validation_passed=False,
1855:         validation_details={"error": "No areas", "type": "empty"}
1856:     )
1857:
1858: # Apply cluster weights
1859: resolved_weights = weights or self._resolve_cluster_weights(cluster_id, cluster_area_scores)
1860: try:
1861:     weighted_score = self.apply_cluster_weights(cluster_area_scores, resolved_weights)
1862:     validation_details["weights"] = {
1863:         "valid": True,
1864:         "weights": resolved_weights if resolved_weights else "equal",
1865:         "score": weighted_score
1866:     }
1867: except WeightValidationError as e:
1868:     logger.error(f"Cluster weight validation failed: {e}")
1869:     return ClusterScore(
1870:         cluster_id=cluster_id,
1871:         cluster_name=cluster_name,
1872:         areas=expected_areas,
1873:         score=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1670_22", 0.0),
1874:         coherence=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1671_26", 0.0),
1875:         variance=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1672_25", 0.0),
1876:         weakest_area=None,
1877:         area_scores=cluster_area_scores,
1878:         validation_passed=False,
1879:         validation_details={"error": str(e), "type": "weights"}
1880:     )
1881:
```

```
1882:     # Analyze coherence and variance metrics
1883:     coherence = self.analyze_coherence(cluster_area_scores)
1884:     scores_array = [a.score for a in cluster_area_scores]
1885:     if scores_array:
1886:         mean_score = sum(scores_array) / len(scores_array)
1887:         variance = sum((score - mean_score) ** 2 for score in scores_array) / len(scores_array)
1888:     else:
1889:         variance = ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "variance", 0.0) # Refactored
1890:     weakest_area = min(cluster_area_scores, key=lambda a: a.score, default=None)
1891:
1892:     std_dev = variance ** ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1689_30", 0.5)
1893:     normalized_std = min(std_dev / self.MAX_SCORE, ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1690_55", 1.0)) if std_dev > 0 else ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1690_80", 0.0)
1894:     penalty_factor = ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1691_25", 1.0) - (normalized_std * self.PENALTY_WEIGHT)
1895:     adjusted_score = weighted_score * penalty_factor
1896:
1897:     validation_details["coherence"] = {
1898:         "value": coherence,
1899:         "interpretation": "high" if coherence > ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1696_52", 0.8) else "medium" if coherence > ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1696_85", 0.6) else "low"
1900:     }
1901:     validation_details["variance"] = variance
1902:     if weakest_area:
1903:         validation_details["weakest_area"] = weakest_area.area_id
1904:     validation_details["imbalance_penalty"] = {
1905:         "std_dev": std_dev,
1906:         "penalty_factor": penalty_factor,
1907:         "raw_score": weighted_score,
1908:         "adjusted_score": adjusted_score,
1909:     }
1910:
1911:     logger.info(
1912:         f"\u234\u223 Cluster {cluster_id} ({cluster_name}): "
1913:         f"score={adjusted_score:.4f}, coherence={coherence:.4f}"
1914:     )
1915:
1916:     return ClusterScore(
1917:         cluster_id=cluster_id,
1918:         cluster_name=cluster_name,
1919:         areas=expected_areas,
1920:         score=adjusted_score,
1921:         coherence=coherence,
1922:         variance=variance,
1923:         weakest_area=weakest_area.area_id if weakest_area else None,
1924:         area_scores=cluster_area_scores,
1925:         validation_passed=True,
1926:         validation_details=validation_details
1927:     )
1928:
1929:     def run(
1930:         self,
1931:         area_scores: list[AreaScore],
1932:         cluster_definitions: list[dict[str, Any]]
```

```
1933:     ) -> list[ClusterScore]:
1934:         """
1935:             Run the cluster aggregation process.
1936:
1937:             Args:
1938:                 area_scores: List of all area scores.
1939:                 cluster_definitions: List of cluster definitions from the monolith.
1940:
1941:             Returns:
1942:                 A list of ClusterScore objects.
1943:             """
1944:             # Create a mapping from area_id to cluster_id
1945:             area_to_cluster = {}
1946:             for cluster in cluster_definitions:
1947:                 for area_id in cluster["policy_area_ids"]:
1948:                     area_to_cluster[area_id] = cluster["cluster_id"]
1949:
1950:             # Assign cluster_id to each area score
1951:             for score in area_scores:
1952:                 score.cluster_id = area_to_cluster.get(score.area_id)
1953:
1954:             def key_func(area_score: AreaScore) -> tuple:
1955:                 return tuple(getattr(area_score, key) for key in self.cluster_group_by_keys)
1956:
1957:             grouped_scores = group_by([s for s in area_scores if hasattr(s, 'cluster_id')], key_func)
1958:
1959:             cluster_scores = []
1960:             for group_key, scores in grouped_scores.items():
1961:                 group_by_values = dict(zip(self.cluster_group_by_keys, group_key, strict=False))
1962:                 score = self.aggregate_cluster(scores, group_by_values)
1963:                 cluster_scores.append(score)
1964:
1965:             return cluster_scores
1966:
1967:     def _resolve_cluster_weights(
1968:         self,
1969:         cluster_id: str,
1970:         area_scores: list[AreaScore],
1971:     ) -> list[float] | None:
1972:         mapping = self.aggregation_settings.cluster_policy_area_weights.get(cluster_id)
1973:         if not mapping:
1974:             return None
1975:
1976:         weights: list[float] = []
1977:         for area_score in area_scores:
1978:             weight = mapping.get(area_score.area_id)
1979:             if weight is None:
1980:                 logger.debug(
1981:                     "Missing weight for area %s in cluster %s \u200\223 falling back to equal weights",
1982:                     area_score.area_id,
1983:                     cluster_id,
1984:                 )
1985:                 return None
1986:             weights.append(weight)
1987:
1988:         total = sum(weights)
```

```
1989:         if total <= 0:
1990:             return None
1991:         return [w / total for w in weights]
1992:
1993: class MacroAggregator:
1994:     """
1995:         Performs holistic macro evaluation (Q305).
1996:
1997:     Responsibilities:
1998:         - Aggregate all cluster scores
1999:         - Calculate cross-cutting coherence
2000:         - Identify systemic gaps
2001:         - Assess strategic alignment
2002:     """
2003:
2004:     def __init__(
2005:         self,
2006:         monolith: dict[str, Any] | None = None,
2007:         abort_on_insufficient: bool = True,
2008:         aggregation_settings: AggregationSettings | None = None,
2009:     ) -> None:
2010:         """
2011:             Initialize macro aggregator.
2012:
2013:         Args:
2014:             monolith: Questionnaire monolith configuration (optional, required for run())
2015:             abort_on_insufficient: Whether to abort on insufficient coverage
2016:
2017:         Raises:
2018:             ValueError: If monolith is None and required for operations
2019:         """
2020:         self.monolith = monolith
2021:         self.abort_on_insufficient = abort_on_insufficient
2022:         self.aggregation_settings = aggregation_settings or AggregationSettings.from_monolith(monolith)
2023:
2024:         # Extract configuration if monolith provided
2025:         if monolith is not None:
2026:             self.scoring_config = monolith["blocks"]["scoring"]
2027:             self.niveles = monolith["blocks"]["niveles_abstraccion"]
2028:         else:
2029:             self.scoring_config = None
2030:             self.niveles = None
2031:
2032:         logger.info("MacroAggregator initialized")
2033:
2034:     def calculate_cross_cutting_coherence(
2035:         self,
2036:         cluster_scores: list[ClusterScore]
2037:     ) -> float:
2038:         """
2039:             Calculate cross-cutting coherence across all clusters.
2040:
2041:         Args:
2042:             cluster_scores: List of cluster scores
2043:
2044:         Returns:
```

```
2045:         Cross-cutting coherence value (0-1)
2046:         """
2047:         scores = [c.score for c in cluster_scores]
2048:
2049:         if len(scores) <= 1:
2050:             return ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1847_19", 1.0)
2051:
2052:         # Calculate mean
2053:         mean = sum(scores) / len(scores)
2054:
2055:         # Calculate standard deviation
2056:         variance = sum((s - mean) ** 2 for s in scores) / len(scores)
2057:         std_dev = variance ** ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1854_30", 0.5)
2058:
2059:         # Convert to coherence
2060:         max_std = 3.0
2061:         coherence = max(ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1858_24", 0.0), ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1858_29", 1.0) - (std_dev / max_std))
2062:
2063:         logger.debug(
2064:             f"Cross-cutting coherence: mean={mean:.4f}, "
2065:             f"std_dev={std_dev:.4f}, coherence={coherence:.4f}"
2066:         )
2067:
2068:     return coherence
2069:
2070:     def identify_systemic_gaps(
2071:         self,
2072:         area_scores: list[AreaScore]
2073:     ) -> list[str]:
2074:         """
2075:             Identify systemic gaps (areas with INSUFICIENTE quality).
2076:
2077:             Args:
2078:                 area_scores: List of area scores
2079:
2080:             Returns:
2081:                 List of area names with systemic gaps
2082:             """
2083:         gaps = []
2084:         for area in area_scores:
2085:             if area.quality_level == "INSUFICIENTE":
2086:                 gaps.append(area.area_name)
2087:                 logger.warning(f"Systemic gap identified: {area.area_name}")
2088:
2089:             logger.info(f"Systemic gaps identified: {len(gaps)}")
2090:         return gaps
2091:
2092:     def assess_strategic_alignment(
2093:         self,
2094:         cluster_scores: list[ClusterScore],
2095:         dimension_scores: list[DimensionScore]
2096:     ) -> float:
2097:         """
2098:             Assess strategic alignment across all levels.
2099:
```

```

2100:     Args:
2101:         cluster_scores: List of cluster scores
2102:         dimension_scores: List of dimension scores
2103:
2104:     Returns:
2105:         Strategic alignment score (0-1)
2106:         """
2107:     # Calculate average cluster coherence
2108:     cluster_coherence = (
2109:         sum(c.coherence for c in cluster_scores) / len(cluster_scores)
2110:         if cluster_scores else ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1907_35", 0.
0)
2111:     )
2112:
2113:     # Calculate dimension validation rate
2114:     validated_dims = sum(1 for d in dimension_scores if d.validation_passed)
2115:     validation_rate = validated_dims / len(dimension_scores) if dimension_scores else ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterA
ggregator.analyze_coherence", "auto_param_L1912_90", 0.0)
2116:
2117:     # Strategic alignment is weighted combination
2118:     alignment = (ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1915_21", 0.6) * cluster_c
oherence) + (ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1915_49", 0.4) * validation_rate)
2119:
2120:     logger.debug(
2121:         f"Strategic alignment: cluster_coherence={cluster_coherence:.4f}, "
2122:         f"validation_rate={validation_rate:.4f}, alignment={alignment:.4f}"
2123:     )
2124:
2125:     return alignment
2126:
2127: def apply_rubric_thresholds(
2128:     self,
2129:     score: float,
2130:     thresholds: dict[str, float] | None = None
2131: ) -> str:
2132:     """
2133:     Apply macro-level rubric thresholds.
2134:
2135:     Args:
2136:         score: Aggregated macro score (0-3 range)
2137:         thresholds: Optional threshold definitions (dict with keys: EXCELENTE, BUENO, ACEPTABLE)
2138:             Each value should be a normalized threshold (0-1 range)
2139:
2140:     Returns:
2141:         Quality level (EXCELENTE, BUENO, ACEPTABLE, INSUFICIENTE)
2142:         """
2143:     # Clamp score to valid range [0, 3]
2144:     clamped_score = max(ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1941_28", 0.0), min
(3.0, score))
2145:
2146:     # Normalize to 0-1 range
2147:     normalized_score = clamped_score / 3.0
2148:
2149:     # Use provided thresholds or defaults
2150:     if thresholds:
2151:         excellent_threshold = thresholds.get('EXCELENTE', ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence")

```

```
, "auto_param_L1948_62", 0.85))
2152:         good_threshold = thresholds.get('BUENO', ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_p
aram_L1949_53", 0.70))
2153:         acceptable_threshold = thresholds.get('ACCEPTABLE', ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence
", "auto_param_L1950_63", 0.55))
2154:     else:
2155:         excellent_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "excellent_threshold", 0.8
5) # Refactored
2156:         good_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "good_threshold", 0.7) # Refact
ored
2157:         acceptable_threshold = ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "acceptable_threshold", 0
.55) # Refactored
2158:
2159:     # Apply thresholds
2160:     if normalized_score >= excellent_threshold:
2161:         quality = "EXCELENTE"
2162:     elif normalized_score >= good_threshold:
2163:         quality = "BUENO"
2164:     elif normalized_score >= acceptable_threshold:
2165:         quality = "ACCEPTABLE"
2166:     else:
2167:         quality = "INSUFICIENTE"
2168:
2169:     logger.debug(
2170:         f"Macro rubric applied: score={score:.4f}, "
2171:         f"normalized={normalized_score:.4f}, quality={quality}"
2172:     )
2173:
2174:     return quality
2175:
2176: def evaluate_macro(
2177:     self,
2178:     cluster_scores: list[ClusterScore],
2179:     area_scores: list[AreaScore],
2180:     dimension_scores: list[DimensionScore]
2181: ) -> MacroScore:
2182:     """
2183:         Perform holistic macro evaluation (Q305).
2184:
2185:     Args:
2186:         cluster_scores: List of cluster scores (MESO level)
2187:         area_scores: List of area scores
2188:         dimension_scores: List of dimension scores
2189:
2190:     Returns:
2191:         MacroScore with holistic evaluation
2192:     """
2193:     logger.info("Performing macro holistic evaluation (Q305)")
2194:
2195:     validation_details = {}
2196:
2197:     if not cluster_scores:
2198:         logger.error("No cluster scores available for macro evaluation")
2199:         return MacroScore(
2200:             score=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1997_22", 0.0),
2201:             quality_level="INSUFICIENTE",
```

```
2202:         cross_cutting_coherence=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L1999_40
", 0.0),
2203:         systemic_gaps=[],
2204:         strategic_alignment=ParameterLoaderV2.get("farfan_core.processing.aggregation.ClusterAggregator.analyze_coherence", "auto_param_L2001_36", 0
.0),
2205:         cluster_scores=[],
2206:         validation_passed=False,
2207:         validation_details={"error": "No clusters", "type": "empty"}
2208:     )
2209:
2210:     # Calculate cross-cutting coherence
2211:     cross_cutting_coherence = self.calculate_cross_cutting_coherence(cluster_scores)
2212:     validation_details["coherence"] = {
2213:         "value": cross_cutting_coherence,
2214:         "clusters": len(cluster_scores)
2215:     }
2216:
2217:     # Identify systemic gaps
2218:     systemic_gaps = self.identify_systemic_gaps(area_scores)
2219:     validation_details["gaps"] = {
2220:         "count": len(systemic_gaps),
2221:         "areas": systemic_gaps
2222:     }
2223:
2224:     # Assess strategic alignment
2225:     strategic_alignment = self.assess_strategic_alignment(
2226:         cluster_scores,
2227:         dimension_scores
2228:     )
2229:     validation_details["alignment"] = {
2230:         "value": strategic_alignment
2231:     }
2232:
2233:     # Calculate overall macro score (weighted average of clusters)
2234:     macro_score = self._calculate_macro_score(cluster_scores)
2235:
2236:     # Apply quality rubric
2237:     quality_level = self.apply_rubric_thresholds(macro_score)
2238:     validation_details["rubric"] = {
2239:         "score": macro_score,
2240:         "quality_level": quality_level
2241:     }
2242:
2243:     logger.info(
2244:         f"\u234\u223 Macro evaluation (Q305): score={macro_score:.4f}, "
2245:         f"quality={quality_level}, coherence={cross_cutting_coherence:.4f}, "
2246:         f"alignment={strategic_alignment:.4f}, gaps={len(systemic_gaps)}"
2247:     )
2248:
2249:     return MacroScore(
2250:         score=macro_score,
2251:         quality_level=quality_level,
2252:         cross_cutting_coherence=cross_cutting_coherence,
2253:         systemic_gaps=systemic_gaps,
2254:         strategic_alignment=strategic_alignment,
2255:         cluster_scores=cluster_scores,
```

```
2256:         validation_passed=True,
2257:         validation_details=validation_details
2258:     )
2259:
2260:     @calibrated_method("farfan_core.processing.aggregation.MacroAggregator._calculate_macro_score")
2261:     def _calculate_macro_score(self, cluster_scores: list[ClusterScore]) -> float:
2262:         weights = self.aggregation_settings.macro_cluster_weights
2263:         if not cluster_scores:
2264:             return ParameterLoaderV2.get("farfan_core.processing.aggregation.MacroAggregator._calculate_macro_score", "auto_param_L2061_19", 0.0)
2265:         if not weights:
2266:             return sum(c.score for c in cluster_scores) / len(cluster_scores)
2267:
2268:         resolved_weights: list[float] = []
2269:         for cluster in cluster_scores:
2270:             weight = weights.get(cluster.cluster_id)
2271:             if weight is None:
2272:                 logger.debug(
2273:                     "Missing macro weight for cluster %s \u200d223 falling back to equal weights",
2274:                     cluster.cluster_id,
2275:                 )
2276:             return sum(c.score for c in cluster_scores) / len(cluster_scores)
2277:         resolved_weights.append(weight)
2278:
2279:         total = sum(resolved_weights)
2280:         if total <= 0:
2281:             return sum(c.score for c in cluster_scores) / len(cluster_scores)
2282:
2283:         normalized = [w / total for w in resolved_weights]
2284:         return sum(
2285:             cluster.score * weight
2286:             for cluster, weight in zip(cluster_scores, normalized, strict=False)
2287:         )
2288:
2289:
2290:
2291: =====
2292: FILE: src/farfan_pipeline/processing/aggregation_provenance.py
2293: =====
2294:
2295: """
2296: Aggregation Provenance System - DAG-based Lineage Tracking
2297:
2298: This module provides full provenance tracking for the hierarchical aggregation pipeline.
2299: It implements W3C PROV-compliant directed acyclic graphs (DAGs) to capture:
2300: - Data lineage: Which micro-questions contributed to which macro-scores
2301: - Operation tracking: How aggregation operations transformed data
2302: - Sensitivity analysis: Which inputs have highest impact on outputs
2303: - Counterfactual reasoning: What-if analysis for policy decisions
2304:
2305: Architecture:
2306: - ProvenanceNode: Immutable record of a score at any aggregation level
2307: - AggregationDAG: NetworkX-based graph with attribution methods
2308: - SHAPAttribution: Shapley value computation for feature importance
2309:
2310: References:
2311: - W3C PROV: https://www.w3.org/TR/prov-overview/
```

```
2312: - Shapley values: Shapley, L.S. (1953). "A value for n-person games"
2313: - NetworkX: Hagberg et al. (2008). "Exploring network structure, dynamics, and function"
2314: """
2315:
2316: from __future__ import annotations
2317:
2318: import hashlib
2319: import json
2320: import logging
2321: from collections import defaultdict
2322: from dataclasses import asdict, dataclass, field
2323: from datetime import datetime, timezone
2324: from typing import Any
2325:
2326: import networkx as nx
2327: import numpy as np
2328:
2329: logger = logging.getLogger(__name__)
2330:
2331:
2332: @dataclass(frozen=True)
2333: class ProvenanceNode:
2334:     """
2335:         Immutable provenance record for a score node in the aggregation DAG.
2336:
2337:     Attributes:
2338:         node_id: Unique identifier (e.g., "Q001", "DIM01_PA05", "CLUSTER_MESO_1")
2339:         level: Abstraction level in hierarchy
2340:         score: Numeric score value
2341:         quality_level: Quality classification (EXCELENTE, BUENO, etc.)
2342:         timestamp: ISO timestamp of computation
2343:         metadata: Additional context (weights, confidence, etc.)
2344:     """
2345:     node_id: str
2346:     level: str # "micro", "dimension", "area", "cluster", "macro"
2347:     score: float
2348:     quality_level: str
2349:     timestamp: str = field(default_factory=lambda: datetime.now(timezone.utc).isoformat())
2350:     metadata: dict[str, Any] = field(default_factory=dict)
2351:
2352:     def to_dict(self) -> dict[str, Any]:
2353:         """Convert to dictionary for serialization."""
2354:         return asdict(self)
2355:
2356:     def compute_hash(self) -> str:
2357:         """Compute deterministic hash for reproducibility."""
2358:         content = json.dumps(
2359:             {
2360:                 "node_id": self.node_id,
2361:                 "level": self.level,
2362:                 "score": self.score,
2363:                 "quality_level": self.quality_level,
2364:             },
2365:             sort_keys=True,
2366:         )
2367:         return hashlib.sha256(content.encode("utf-8")).hexdigest()[:16]
```

```
2368:  
2369:  
2370: @dataclass  
2371: class AggregationEdge:  
2372:     """  
2373:         Edge in the provenance DAG representing an aggregation operation.  
2374:  
2375:         Attributes:  
2376:             source_id: Parent node contributing to aggregation  
2377:             target_id: Child node receiving aggregated score  
2378:             operation: Type of aggregation (weighted_average, choquet, max, etc.)  
2379:             weight: Contribution weight (for weighted operations)  
2380:             timestamp: When this edge was created  
2381:         """  
2382:         source_id: str  
2383:         target_id: str  
2384:         operation: str  
2385:         weight: float  
2386:         timestamp: str = field(default_factory=lambda: datetime.now(timezone.utc).isoformat())  
2387:         metadata: dict[str, Any] = field(default_factory=dict)  
2388:  
2389:  
2390: class AggregationDAG:  
2391:     """  
2392:         Directed Acyclic Graph for full provenance tracking of aggregation pipeline.  
2393:  
2394:         This class maintains the complete lineage of how micro-question scores  
2395:         propagate through dimension \206\222 area \206\222 cluster \206\222 macro aggregation.  
2396:  
2397:         Features:  
2398:             - Cycle detection (enforces DAG property)  
2399:             - Topological sorting for dependency resolution  
2400:             - Shapley value attribution for sensitivity analysis  
2401:             - GraphML export for visualization in Gephi/Cytoscape  
2402:             - W3C PROV-compliant serialization  
2403:         """  
2404:  
2405:     def __init__(self):  
2406:         """Initialize empty DAG."""  
2407:         self.graph: nx.DiGraph = nx.DiGraph()  
2408:         self.nodes: dict[str, ProvenanceNode] = {}  
2409:         self.edges: list[AggregationEdge] = []  
2410:         logger.info("AggregationDAG initialized")  
2411:  
2412:     def add_node(self, node: ProvenanceNode) -> None:  
2413:         """  
2414:             Add a provenance node to the DAG.  
2415:  
2416:             Args:  
2417:                 node: ProvenanceNode to add  
2418:  
2419:             Raises:  
2420:                 ValueError: If node_id already exists  
2421:         """  
2422:         if node.node_id in self.nodes:  
2423:             logger.warning(f"Node {node.node_id} already exists, skipping")
```

```
2424:         return
2425:
2426:     self.nodes[node.node_id] = node
2427:     self.graph.add_node(
2428:         node.node_id,
2429:         level=node.level,
2430:         score=node.score,
2431:         quality=node.quality_level,
2432:         timestamp=node.timestamp,
2433:         metadata=node.metadata,
2434:     )
2435:     logger.debug(f"Added node {node.node_id} (level={node.level}, score={node.score:.2f})")
2436:
2437: def add_aggregation_edge(
2438:     self,
2439:     source_ids: list[str],
2440:     target_id: str,
2441:     operation: str,
2442:     weights: list[float],
2443:     metadata: dict[str, Any] | None = None,
2444: ) -> None:
2445:     """
2446:     Record an aggregation operation: sources \206\222 target.
2447:
2448:     Args:
2449:         source_ids: List of source node IDs (e.g., micro-questions)
2450:         target_id: Target node ID (e.g., dimension score)
2451:         operation: Aggregation type (weighted_average, choquet, etc.)
2452:         weights: Contribution weights for each source
2453:         metadata: Additional operation metadata
2454:
2455:     Raises:
2456:         ValueError: If weights don't match sources or cycle detected
2457:     """
2458:     if len(source_ids) != len(weights):
2459:         raise ValueError(
2460:             f"Mismatch: {len(source_ids)} sources but {len(weights)} weights"
2461:         )
2462:
2463:     if metadata is None:
2464:         metadata = {}
2465:
2466:     for source_id, weight in zip(source_ids, weights):
2467:         if source_id not in self.graph:
2468:             logger.warning(f"Source node {source_id} not found, adding placeholder")
2469:             self.graph.add_node(source_id, level="unknown", score=0.0)
2470:
2471:         if target_id not in self.graph:
2472:             logger.warning(f"Target node {target_id} not found, adding placeholder")
2473:             self.graph.add_node(target_id, level="unknown", score=0.0)
2474:
2475:         edge = AggregationEdge(
2476:             source_id=source_id,
2477:             target_id=target_id,
2478:             operation=operation,
2479:             weight=weight,
```

```
2480:             metadata=metadata,
2481:         )
2482:         self.edges.append(edge)
2483:
2484:         self.graph.add_edge(
2485:             source_id,
2486:             target_id,
2487:             operation=operation,
2488:             weight=weight,
2489:             timestamp=edge.timestamp,
2490:             metadata=metadata,
2491:         )
2492:
2493:     # Verify DAG property (no cycles)
2494:     if not nx.is_directed_acyclic_graph(self.graph):
2495:         # Rollback last edges
2496:         for source_id in source_ids:
2497:             if self.graph.has_edge(source_id, target_id):
2498:                 self.graph.remove_edge(source_id, target_id)
2499:         raise ValueError(f"Adding edges {source_ids} \u2013 {target_id} would create a cycle")
2500:
2501:     logger.info(
2502:         f"Added aggregation: {len(source_ids)} sources \u2013 {target_id} "
2503:         f"(operation={operation})"
2504:     )
2505:
2506:     def trace_lineage(self, target_id: str) -> dict[str, Any]:
2507:         """
2508:             Trace complete lineage of a target node.
2509:
2510:             Returns all ancestor nodes, the aggregation path, and sensitivity metrics.
2511:
2512:             Args:
2513:                 target_id: Node ID to trace
2514:
2515:             Returns:
2516:                 Dictionary with:
2517:                     - ancestors: Set of all ancestor node IDs
2518:                     - path: Topologically sorted path from sources to target
2519:                     - depth: Maximum path length from any micro-question
2520:                     - breadth: Number of unique micro-questions contributing
2521:
2522:             Raises:
2523:                 ValueError: If target_id not in graph
2524:         """
2525:         if target_id not in self.graph:
2526:             raise ValueError(f"Node {target_id} not found in DAG")
2527:
2528:         ancestors = nx.ancestors(self.graph, target_id)
2529:         subgraph = self.graph.subgraph(ancestors | {target_id})
2530:
2531:         # Compute metrics
2532:         depth = nx.dag_longest_path_length(subgraph) if subgraph.nodes else 0
2533:
2534:         # Count micro-questions (level="micro")
2535:         micro_nodes = [
```

```
2536:         n for n in ancestors
2537:         if self.graph.nodes[n].get("level") == "micro"
2538:     ]
2539:
2540:     # Get topological path
2541:     topo_path = list(nx.topological_sort(subgraph))
2542:
2543:     return {
2544:         "target_id": target_id,
2545:         "ancestor_count": len(ancestors),
2546:         "ancestors": sorted(ancestors),
2547:         "topological_path": topo_path,
2548:         "depth": depth,
2549:         "micro_question_count": len(micro_nodes),
2550:         "micro_questions": sorted(micro_nodes),
2551:     }
2552:
2553:     def compute_shapley_attribution(self, target_id: str) -> dict[str, float]:
2554:         """
2555:             Compute Shapley values for feature attribution.
2556:
2557:             Shapley values represent the marginal contribution of each source node
2558:             to the target score, accounting for all possible coalitions.
2559:
2560:             Args:
2561:                 target_id: Node to attribute
2562:
2563:             Returns:
2564:                 Dictionary mapping source node IDs to Shapley values (sum = target score)
2565:
2566:             Note:
2567:                 This is an exact computation for weighted averages. For non-linear
2568:                 aggregations (Choquet), this uses kernel SHAP approximation.
2569:         """
2570:         if target_id not in self.graph:
2571:             raise ValueError(f"Node {target_id} not found in DAG")
2572:
2573:         # Get direct predecessors (sources)
2574:         predecessors = list(self.graph.predecessors(target_id))
2575:
2576:         if not predecessors:
2577:             logger.warning(f"Node {target_id} has no predecessors")
2578:             return {}
2579:
2580:         # Extract weights and scores
2581:         weights = []
2582:         scores = []
2583:         for pred in predecessors:
2584:             edge_data = self.graph.get_edge_data(pred, target_id)
2585:             weights.append(edge_data.get("weight", 0.0))
2586:             scores.append(self.graph.nodes[pred].get("score", 0.0))
2587:
2588:         weights = np.array(weights)
2589:         scores = np.array(scores)
2590:
2591:         # For weighted average, Shapley value = weight / 2^7 score
```

```
2592:     # This is exact because weighted average is a linear function
2593:     shapley_values = weights * scores
2594:
2595:     # Normalize to sum to target score
2596:     target_score = self.graph.nodes[target_id].get("score", 0.0)
2597:     if np.sum(shapley_values) > 0:
2598:         shapley_values = shapley_values * (target_score / np.sum(shapley_values))
2599:
2600:     attribution = {
2601:         pred: float(shap_val)
2602:         for pred, shap_val in zip(predecessors, shapley_values)
2603:     }
2604:
2605:     logger.debug(
2606:         f"Shapley attribution for {target_id}: "
2607:         f"{len(attribution)} sources, sum={sum(attribution.values()):.4f}"
2608:     )
2609:
2610:     return attribution
2611:
2612: def get_critical_path(self, target_id: str, top_k: int = 5) -> list[tuple[str, float]]:
2613:     """
2614:     Identify the most critical source nodes for a target.
2615:
2616:     Uses Shapley values to rank sources by importance.
2617:
2618:     Args:
2619:         target_id: Target node
2620:         top_k: Number of top sources to return
2621:
2622:     Returns:
2623:         List of (node_id, shapley_value) tuples, sorted by importance
2624:     """
2625:     attribution = self.compute_shapley_attribution(target_id)
2626:
2627:     # Sort by absolute Shapley value (handle negative contributions)
2628:     sorted_attribution = sorted(
2629:         attribution.items(),
2630:         key=lambda x: abs(x[1]),
2631:         reverse=True,
2632:     )
2633:
2634:     return sorted_attribution[:top_k]
2635:
2636: def export_graphml(self, path: str) -> None:
2637:     """
2638:     Export DAG to GraphML format for visualization.
2639:
2640:     Compatible with:
2641:     - Gephi: https://gephi.org/
2642:     - Cytoscape: https://cytoscape.org/
2643:     - yEd: https://www.yworks.com/products/yed
2644:
2645:     Args:
2646:         path: Output file path (e.g., "aggregation_dag.graphml")
2647:     """
```

```

2648:     # Create a copy with serialized dict attributes (GraphML doesn't support dicts)
2649:     g_copy = self.graph.copy()
2650:     for node, data in g_copy.nodes(data=True):
2651:         if "metadata" in data and isinstance(data["metadata"], dict):
2652:             data["metadata_json"] = json.dumps(data["metadata"])
2653:             del data["metadata"]
2654:     for u, v, data in g_copy.edges(data=True):
2655:         if "metadata" in data and isinstance(data["metadata"], dict):
2656:             data["metadata_json"] = json.dumps(data["metadata"])
2657:             del data["metadata"]
2658:         if "weights" in data and isinstance(data["weights"], list):
2659:             data["weights_json"] = json.dumps(data["weights"])
2660:             del data["weights"]
2661:     nx.write_graphml(g_copy, path)
2662:     logger.info(
2663:         f"Exported DAG to {path}: "
2664:         f"{self.graph.number_of_nodes()} nodes, "
2665:         f"{self.graph.number_of_edges()} edges"
2666:     )
2667:
2668: def export_prov_json(self, path: str) -> None:
2669:     """
2670:     Export to W3C PROV-JSON format.
2671:
2672:     Spec: https://www.w3.org/Submission/prov-json/
2673:
2674:     Args:
2675:         path: Output file path (e.g., "provenance.json")
2676:     """
2677:     prov_doc = {
2678:         "prefix": {
2679:             "prov": "http://www.w3.org/ns/prov#",
2680:             "farfan": "http://farfan.org/ns/aggregation#",
2681:         },
2682:         "entity": {},
2683:         "activity": {},
2684:         "wasGeneratedBy": {},
2685:         "used": {}
2686:     }
2687:
2688:     # Entities: All nodes
2689:     for node_id, node_data in self.graph.nodes(data=True):
2690:         prov_doc["entity"][node_id] = {
2691:             "prov:type": "farfan:ScoreEntity",
2692:             "farfan:level": node_data.get("level"),
2693:             "farfan:score": node_data.get("score"),
2694:             "farfan:quality": node_data.get("quality"),
2695:             "prov:generatedAtTime": node_data.get("timestamp"),
2696:         }
2697:
2698:     # Activities: Aggregation operations
2699:     for idx, edge in enumerate(self.edges):
2700:         activity_id = f"agg_{idx}"
2701:         prov_doc["activity"][activity_id] = {
2702:             "prov:type": f"farfan:{edge.operation}",
2703:             "farfan:weight": edge.weight,

```

```
2704:         "prov:startedAtTime": edge.timestamp,
2705:     }
2706:
2707:     # wasGeneratedBy: target was generated by this activity
2708:     if edge.target_id not in prov_doc["wasGeneratedBy"]:
2709:         prov_doc["wasGeneratedBy"][edge.target_id] = []
2710:         prov_doc["wasGeneratedBy"][edge.target_id].append(activity_id)
2711:
2712:     # used: activity used source
2713:     if activity_id not in prov_doc["used"]:
2714:         prov_doc["used"][activity_id] = []
2715:         prov_doc["used"][activity_id].append(edge.source_id)
2716:
2717:     with open(path, "w", encoding="utf-8") as f:
2718:         json.dump(prov_doc, f, indent=2, ensure_ascii=False)
2719:
2720:     logger.info(f"Exported PROV-JSON to {path}")
2721:
2722: def get_statistics(self) -> dict[str, Any]:
2723:     """
2724:     Get DAG statistics for monitoring and validation.
2725:
2726:     Returns:
2727:         Dictionary with graph metrics
2728:     """
2729:     return {
2730:         "node_count": self.graph.number_of_nodes(),
2731:         "edge_count": self.graph.number_of_edges(),
2732:         "max_depth": nx.dag_longest_path_length(self.graph) if self.graph.nodes else 0,
2733:         "is_dag": nx.is_directed_acyclic_graph(self.graph),
2734:         "weakly_connected_components": nx.number_weakly_connected_components(self.graph),
2735:         "nodes_by_level": self._count_by_level(),
2736:     }
2737:
2738:     def _count_by_level(self) -> dict[str, int]:
2739:         """Count nodes by abstraction level."""
2740:         counts = defaultdict(int)
2741:         for node_id in self.graph.nodes:
2742:             level = self.graph.nodes[node_id].get("level", "unknown")
2743:             counts[level] += 1
2744:         return dict(counts)
2745:
2746:
2747:     __all__ = [
2748:         "ProvenanceNode",
2749:         "AggregationEdge",
2750:         "AggregationDAG",
2751:     ]
2752:
2753:
2754:
2755: =====
2756: FILE: src/farfan_pipeline/processing/choquet_adapter.py
2757: =====
2758:
2759: """
```

```
2760: Choquet Integral Adapter for Processing Aggregation
2761:
2762: This module adapts the calibration ChoquetAggregator for use in processing-level
2763: score aggregation. It handles the impedance mismatch between:
2764: - Calibration layer scores (8 layers: @b, @chain, @q, @d, @p, @C, @u, @m)
2765: - Processing dimension/area scores (numeric scores with weights)
2766:
2767: Key adaptations:
2768: - Maps dimension/question indices to pseudo-layer IDs
2769: - Converts score lists + weights to layer score dictionaries
2770: - Preserves non-linear interaction semantics from Choquet integral
2771:
2772: References:
2773: - Grabisch, M. (2016). "Set Functions, Games and Capacities in Decision Making"
2774: - Choquet, G. (1953). "Theory of capacities"
2775: """
2776:
2777: from __future__ import annotations
2778:
2779: import logging
2780: from typing import Any
2781:
2782: import numpy as np
2783:
2784: logger = logging.getLogger(__name__)
2785:
2786:
2787: class ChoquetProcessingAdapter:
2788:     """
2789:         Adapter for using Choquet integral in processing aggregation.
2790:
2791:     This class bridges the gap between calibration (layer-based) and
2792:     processing (score list) contexts.
2793: """
2794:
2795:     def __init__(
2796:         self,
2797:         linear_weights: dict[str, float] | None = None,
2798:         interaction_weights: dict[tuple[str, str], float] | None = None,
2799:     ):
2800:         """
2801:             Initialize Choquet adapter.
2802:
2803:             Args:
2804:                 linear_weights: Weights for linear terms (default: uniform)
2805:                 interaction_weights: Weights for interaction terms (default: none)
2806: """
2807:         self.linear_weights = linear_weights or {}
2808:         self.interaction_weights = interaction_weights or {}
2809:         logger.info(
2810:             f"ChoquetProcessingAdapter initialized: "
2811:             f"{len(self.linear_weights)} linear, "
2812:             f"{len(self.interaction_weights)} interaction terms"
2813:         )
2814:
2815:     def aggregate(
```

```
2816:         self,
2817:         scores: list[float],
2818:         weights: list[float] | None = None,
2819:         interaction_pairs: list[tuple[int, int]] | None = None,
2820:     ) -> float:
2821:         """
2822:             Aggregate scores using Choquet 2-additive integral.
2823:
2824:             Formula:
2825:                  $C(x) = \sum w_i * x_i + \sum w_{ij} * \min(x_i, x_j)$ 
2826:
2827:             Where:
2828:                 - First sum: Linear terms (weighted individual scores)
2829:                 - Second sum: Interaction terms (synergy/redundancy between pairs)
2830:
2831:             Args:
2832:                 scores: List of numeric scores
2833:                 weights: Linear weights for each score (default: uniform)
2834:                 interaction_pairs: List of (i, j) index pairs for interactions
2835:
2836:             Returns:
2837:                 Aggregated score
2838:
2839:             Raises:
2840:                 ValueError: If weights mismatch or scores empty
2841:             """
2842:             if not scores:
2843:                 raise ValueError("Cannot aggregate empty score list")
2844:
2845:             n = len(scores)
2846:             scores_arr = np.array(scores, dtype=np.float64)
2847:
2848:             # Default: uniform weights
2849:             if weights is None:
2850:                 weights_arr = np.ones(n) / n
2851:             else:
2852:                 if len(weights) != n:
2853:                     raise ValueError(f"Weight count {len(weights)} != score count {n}")
2854:                 weights_arr = np.array(weights, dtype=np.float64)
2855:                 # Normalize to sum to 1.0
2856:                 weights_arr = weights_arr / np.sum(weights_arr)
2857:
2858:             # STEP 1: Linear contribution
2859:             #  $C_{linear} = \sum w_i * x_i$ 
2860:             linear_contrib = float(np.sum(weights_arr * scores_arr))
2861:
2862:             logger.debug(f"Choquet linear contribution: {linear_contrib:.4f}")
2863:
2864:             # STEP 2: Interaction contribution
2865:             #  $C_{interaction} = \sum w_{ij} * \min(x_i, x_j)$ 
2866:             interaction_contrib = 0.0
2867:
2868:             if interaction_pairs:
2869:                 for i, j in interaction_pairs:
2870:                     if i >= n or j >= n:
2871:                         logger.warning(f"Interaction pair ({i}, {j}) out of bounds (n={n})")
```

```
2872:         continue
2873:
2874:         # Interaction weight (default: small positive for synergy)
2875:         pair_key = (str(i), str(j))
2876:         w_ij = self.interaction_weights.get(pair_key, 0.05)
2877:
2878:         # Choquet uses min() for weakest-link principle
2879:         min_score = min(scores_arr[i], scores_arr[j])
2880:         contribution = w_ij * min_score
2881:         interaction_contrib += contribution
2882:
2883:         logger.debug(
2884:             f"Interaction ({i}, {j}): "
2885:             f"min({scores_arr[i]:.2f}, {scores_arr[j]:.2f}) = {min_score:.2f}, "
2886:             f"w={w_ij:.3f}, contrib={contribution:.4f}"
2887:         )
2888:
2889:     logger.debug(f"Choquet interaction contribution: {interaction_contrib:.4f}")
2890:
2891:     # STEP 3: Total score
2892:     total = linear_contrib + interaction_contrib
2893:
2894:     logger.info(
2895:         f"Choquet aggregation: {n} inputs \u2063\222 {total:.4f} "
2896:         f"(linear={linear_contrib:.4f}, interaction={interaction_contrib:.4f})"
2897:     )
2898:
2899:     return float(total)
2900:
2901: @staticmethod
2902: def detect_interactions(
2903:     scores: list[float],
2904:     weights: list[float],
2905:     correlation_threshold: float = 0.7,
2906) -> list[tuple[int, int]]:
2907: """
2908:     Auto-detect interaction pairs based on score correlation.
2909:
2910:     Pairs with high correlation should be penalized (redundancy),
2911:     pairs with negative correlation should be rewarded (complementarity).
2912:
2913:     Args:
2914:         scores: Score values
2915:         weights: Score weights
2916:         correlation_threshold: Min correlation for interaction
2917:
2918:     Returns:
2919:         List of (i, j) pairs with significant interactions
2920:
2921:     Note:
2922:         This is a heuristic. For production, use domain knowledge to
2923:         specify interactions explicitly.
2924: """
2925:     n = len(scores)
2926:     pairs = []
2927:
```

```
2928:         # For now, use a simple heuristic: top-weighted pairs
2929:         # In practice, you'd analyze historical data for correlations
2930:         sorted_indices = np.argsort(weights) [::-1]    # Descending
2931:
2932:         # Consider top 3 interactions among high-weight inputs
2933:         if n >= 2:
2934:             for i in range(min(3, n - 1)):
2935:                 idx_i = sorted_indices[i]
2936:                 idx_j = sorted_indices[i + 1]
2937:                 pairs.append((int(idx_i), int(idx_j)))
2938:
2939:             logger.debug(f"Auto-detected {len(pairs)} interaction pairs")
2940:         return pairs
2941:
2942:
2943: def create_default_choquet_adapter(n_inputs: int) -> ChoquetProcessingAdapter:
2944:     """
2945:         Create a Choquet adapter with sensible defaults.
2946:
2947:         Default configuration:
2948:             - Uniform linear weights (each input gets 1/n)
2949:             - Small interaction weights (0.05) for adjacent pairs
2950:             - Synergistic interactions (positive weights)
2951:
2952:         Args:
2953:             n_inputs: Number of input scores
2954:
2955:         Returns:
2956:             Configured ChoquetProcessingAdapter
2957:     """
2958:     # Linear weights: uniform
2959:     linear_weights = {str(i): 1.0 / n_inputs for i in range(n_inputs)}
2960:
2961:     # Interaction weights: adjacent pairs only (locality assumption)
2962:     interaction_weights = {}
2963:     for i in range(n_inputs - 1):
2964:         pair_key = (str(i), str(i + 1))
2965:         interaction_weights[pair_key] = 0.05    # Small synergy
2966:
2967:     logger.info(
2968:         f"Created default Choquet adapter: "
2969:         f"{n_inputs} inputs, {len(interaction_weights)} interactions"
2970:     )
2971:
2972:     return ChoquetProcessingAdapter(
2973:         linear_weights=linear_weights,
2974:         interaction_weights=interaction_weights,
2975:     )
2976:
2977:
2978: __all__ = [
2979:     "ChoquetProcessingAdapter",
2980:     "create_default_choquet_adapter",
2981: ]
```

```
2984:  
2985: =====  
2986: FILE: src/farfán_pipeline/processing/converter.py  
2987: =====  
2988:  
2989: """  
2990: SmartChunk to CanonPolicyPackage Converter  
2991: =====  
2992:  
2993: This module provides the critical bridge layer between the SPC (Smart Policy Chunks)  
2994: phase-one output and the CanonPolicyPackage format expected by the orchestrator.  
2995:  
2996: Architecture:  
2997:     SmartPolicyChunk (from StrategicChunkingSystem)  
2998:         à\206\223  
2999:     SmartChunkConverter (this module)  
3000:         à\206\223  
3001:     CanonPolicyPackage (for SPCAdapter and Orchestrator)  
3002:  
3003: Key Responsibilities:  
3004: 1. Convert SmartPolicyChunk dataclass to Chunk dataclass  
3005: 2. Map chunk_type (8 types) to resolution (MICRO/MESO/MACRO)  
3006: 3. Extract policy/time/geo facets from SPC rich data  
3007: 4. Build ChunkGraph with edges from related_chunks  
3008: 5. Preserve SPC rich data in metadata for executor access  
3009: 6. Generate quality metrics and integrity index  
3010: """  
3011:  
3012: from __future__ import annotations  
3013:  
3014: import hashlib  
3015: import json  
3016: import logging  
3017: from typing import TYPE_CHECKING, Any  
3018:  
3019: from farfan_pipeline.core.calibration.decorators import calibrated_method  
3020: from farfan_pipeline.processing.cpp_ingestion.models import (  
3021:     KPI,  
3022:     Budget,  
3023:     CanonPolicyPackage,  
3024:     Chunk,  
3025:     ChunkGraph,  
3026:     ChunkResolution,  
3027:     Confidence,  
3028:     Entity,  
3029:     GeoFacet,  
3030:     IntegrityIndex,  
3031:     PolicyFacet,  
3032:     PolicyManifest,  
3033:     ProvenanceMap,  
3034:     QualityMetrics,  
3035:     TextSpan,  
3036:     TimeFacet,  
3037: )  
3038:  
3039: if TYPE_CHECKING:
```

```
3040:     # Avoid runtime import of SmartPolicyChunk (heavy dependencies)
3041:     from typing import Protocol
3042:
3043:     class SmartPolicyChunkProtocol(Protocol):
3044:         """Protocol for SmartPolicyChunk to avoid circular imports"""
3045:         chunk_id: str
3046:         document_id: str
3047:         content_hash: str
3048:         text: str
3049:         normalized_text: str
3050:         semantic_density: float
3051:         section_hierarchy: list[str]
3052:         document_position: tuple[int, int]
3053:         chunk_type: Any # ChunkType enum
3054:         causal_chain: list[Any]
3055:         policy_entities: list[Any]
3056:         related_chunks: list[tuple[str, float]]
3057:         confidence_metrics: dict[str, float]
3058:         coherence_score: float
3059:         completeness_index: float
3060:         strategic_importance: float
3061:
3062: logger = logging.getLogger(__name__)
3063:
3064:
3065: class SmartChunkConverter:
3066:     """
3067:     Converts SmartPolicyChunk instances to CanonPolicyPackage format.
3068:
3069:     This converter is the critical bridge that enables SPC phase-one output
3070:     to be consumed by the orchestrator and its executors.
3071:     """
3072:
3073:     # Mapping from ChunkType to ChunkResolution
3074:     CHUNK_TYPE_TO_RESOLUTION = {
3075:         'DIAGNOSTICO': ChunkResolution.MESO,
3076:         'ESTRATEGIA': ChunkResolution.MACRO,
3077:         'METRICA': ChunkResolution.MICRO,
3078:         'FINANCIERO': ChunkResolution.MICRO,
3079:         'NORMATIVO': ChunkResolution.MESO,
3080:         'OPERATIVO': ChunkResolution.MICRO,
3081:         'EVALUACION': ChunkResolution.MESO,
3082:         'MIXTO': ChunkResolution.MESO,
3083:     }
3084:
3085:     def __init__(self) -> None:
3086:         """Initialize the converter."""
3087:         self.logger = logging.getLogger(self.__class__.__name__)
3088:
3089:     def convert_to_canon_package(
3090:         self,
3091:         smart_chunks: list[Any], # List[SmartPolicyChunk]
3092:         document_metadata: dict[str, Any]
3093:     ) -> CanonPolicyPackage:
3094:         """
3095:             Convert list of SmartPolicyChunk to CanonPolicyPackage.
```

```
3096:  
3097:     Args:  
3098:         smart_chunks: List of SmartPolicyChunk instances from StrategicChunkingSystem  
3099:         document_metadata: Document-level metadata (id, title, version, etc.)  
3100:  
3101:     Returns:  
3102:         CanonPolicyPackage ready for orchestrator consumption  
3103:  
3104:     Raises:  
3105:         ValueError: If smart_chunks is empty or invalid  
3106:         """  
3107:         # Defensive validation: ensure smart_chunks is non-empty  
3108:         if not smart_chunks or len(smart_chunks) == 0:  
3109:             raise ValueError(  
3110:                 "Cannot convert empty smart_chunks list to CanonPolicyPackage. "  
3111:                 "Minimum 1 chunk required from StrategicChunkingSystem."  
3112:             )  
3113:  
3114:         # Defensive validation: check critical attributes on first chunk  
3115:         first_chunk = smart_chunks[0]  
3116:         requiredAttrs = ['chunk_id', 'document_id', 'text', 'document_position', 'chunk_type']  
3117:         missingAttrs = [attr for attr in requiredAttrs if not hasattr(first_chunk, attr)]  
3118:  
3119:         if missingAttrs:  
3120:             raise ValueError(  
3121:                 f"SmartPolicyChunk missing critical attributes: {missingAttrs}. "  
3122:                 f"Ensure StrategicChunkingSystem produced valid SmartPolicyChunk instances. "  
3123:                 f"Chunk type: {type(first_chunk)}")  
3124:         )  
3125:  
3126:         self.logger.info(f"Converting {len(smart_chunks)} SmartPolicyChunks to CanonPolicyPackage")  
3127:  
3128:         # Build ChunkGraph  
3129:         chunk_graph = ChunkGraph()  
3130:  
3131:         # Convert each SmartPolicyChunk to Chunk  
3132:         chunk_hashes = {}  
3133:         all_axes = set()  
3134:         all_programs = set()  
3135:         all_projects = set()  
3136:         all_years = set()  
3137:         all_territories = set()  
3138:  
3139:         for smart_chunk in smart_chunks:  
3140:             # Convert to Chunk  
3141:             chunk = self._convert_smart_chunk_to_chunk(smart_chunk)  
3142:  
3143:             # Add to ChunkGraph  
3144:             chunk_graph.chunks[chunk.id] = chunk  
3145:             chunk_hashes[chunk.id] = chunk.bytes_hash  
3146:  
3147:             # Collect manifest data  
3148:             all_axes.update(chunk.policy_facets.axes)  
3149:             all_programs.update(chunk.policy_facets.programs)  
3150:             all_projects.update(chunk.policy_facets.projects)  
3151:             all_years.update(chunk.time_facets.years)
```

```

3152:         all_territories.update(chunk.geo_facets.territories)
3153:
3154:     # Build edges from related_chunks
3155:     for smart_chunk in smart_chunks:
3156:         if hasattr(smart_chunk, 'related_chunks') and smart_chunk.related_chunks:
3157:             for related_id, similarity in smart_chunk.related_chunks[:5]: # Top 5
3158:                 # Only add edge if target chunk exists
3159:                 if related_id in chunk_graph.chunks:
3160:                     edge = (smart_chunk.chunk_id, related_id, f"semantic_similarity_{similarity:.2f}")
3161:                     chunk_graph.edges.append(edge)
3162:
3163:             self.logger.info(f"Built ChunkGraph with {len(chunk_graph.chunks)} chunks and {len(chunk_graph.edges)} edges")
3164:
3165:     # Create PolicyManifest
3166:     policy_manifest = PolicyManifest(
3167:         axes=sorted(all_axes),
3168:         programs=sorted(all_programs),
3169:         projects=sorted(all_projects),
3170:         years=sorted(all_years),
3171:         territories=sorted(all_territories),
3172:         indicators=[], # Would extract from KPIs if available
3173:         budget_rows=sum(1 for c in chunk_graph.chunks.values() if c.budget is not None)
3174:     )
3175:
3176:     # Calculate QualityMetrics
3177:     quality_metrics = self._calculate_quality_metrics(smart_chunks, chunk_graph)
3178:
3179:     # Generate IntegrityIndex
3180:     integrity_index = self._generate_integrity_index(chunk_hashes, document_metadata)
3181:
3182:     # Preserve SPC rich data in metadata
3183:     enriched_metadata = self._preserve_spc_rich_data(smart_chunks, document_metadata)
3184:
3185:     # Build CanonPolicyPackage
3186:     canon_package = CanonPolicyPackage(
3187:         schema_version="SPC-2025.1",
3188:         chunk_graph=chunk_graph,
3189:         policy_manifest=policy_manifest,
3190:         quality_metrics=quality_metrics,
3191:         integrity_index=integrity_index,
3192:         metadata=enriched_metadata
3193:     )
3194:
3195:     self.logger.info("Successfully converted to CanonPolicyPackage")
3196:     return canon_package
3197:
3198: @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._convert_smart_chunk_to_chunk")
3199: def _convert_smart_chunk_to_chunk(self, smart_chunk: Any) -> Chunk:
3200:
3201:     """
3201:     Convert a single SmartPolicyChunk to Chunk.
3202:
3203:     Maps fields from SPC rich format to orchestrator-compatible format.
3204:     """
3205:
3206:     # Determine resolution from chunk_type
3207:     chunk_type_str = smart_chunk.chunk_type.value if hasattr(smart_chunk.chunk_type, 'value') else str(smart_chunk.chunk_type)
3207:     resolution = self.CHUNK_TYPE_TO_RESOLUTION.get(chunk_type_str.upper(), ChunkResolution.MESO)

```

```
3208:  
3209:     # Extract policy facets  
3210:     policy_facets = self._extract_policy_facets(smart_chunk)  
3211:  
3212:     # Extract time facets  
3213:     time_facets = self._extract_time_facets(smart_chunk)  
3214:  
3215:     # Extract geo facets  
3216:     geo_facets = self._extract_geo_facets(smart_chunk)  
3217:  
3218:     # Build confidence from SPC metrics  
3219:     confidence = Confidence(  
3220:         layout=ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._convert_smart_chunk_to_chunk", "auto_param_L23  
2_19", 1.0), # SPC doesn't distinguish these  
3221:         ocr=smart_chunk.confidence_metrics.get('extraction_confidence', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunk  
Converter._convert_smart_chunk_to_chunk", "auto_param_L233_76", 0.95)),  
3222:         typing=smart_chunk.coherence_score  
3223:     )  
3224:  
3225:     # Create provenance  
3226:     provenance = self._build_provenance(smart_chunk)  
3227:  
3228:     # Extract entities  
3229:     entities = self._extract_entities(smart_chunk)  
3230:  
3231:     # Extract budget if available  
3232:     budget = self._extract_budget(smart_chunk)  
3233:  
3234:     # Extract KPI if available  
3235:     kpi = self._extract_kpi(smart_chunk)  
3236:  
3237:     # Build Chunk  
3238:     return Chunk(  
3239:         id=smart_chunk.chunk_id,  
3240:         text=smart_chunk.text,  
3241:         text_span=TextSpan(  
3242:             start=smart_chunk.document_position[0],  
3243:             end=smart_chunk.document_position[1]  
3244:         ),  
3245:         resolution=resolution,  
3246:         bytes_hash=smart_chunk.content_hash,  
3247:         policy_area_id=getattr(smart_chunk, 'policy_area_id', None), # PA01-PA10  
3248:         dimension_id=getattr(smart_chunk, 'dimension_id', None), # DIM01-DIM06  
3249:         policy_facets=policy_facets,  
3250:         time_facets=time_facets,  
3251:         geo_facets=geo_facets,  
3252:         confidence=confidence,  
3253:         provenance=provenance,  
3254:         budget=budget,  
3255:         kpi=kpi,  
3256:         entities=entities  
3257:     )  
3258:  
3259: @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_policy_facets")  
3260: def _extract_policy_facets(self, smart_chunk: Any) -> PolicyFacet:  
3261:     """Extract policy facets from SPC strategic_context and section_hierarchy."""
```

```

3262:     axes = []
3263:     programs = []
3264:     projects = []
3265:
3266:     # Extract from strategic_context if available
3267:     if hasattr(smart_chunk, 'strategic_context') and smart_chunk.strategic_context:
3268:         ctx = smart_chunk.strategic_context
3269:         # strategic_context might have policy_intent, implementation_phase
3270:         if hasattr(ctx, 'policy_intent'):
3271:             axes.append(ctx.policy_intent[:50]) # Truncate if too long
3272:         if hasattr(ctx, 'implementation_phase'):
3273:             programs.append(ctx.implementation_phase[:50])
3274:
3275:     # Extract from section_hierarchy
3276:     if hasattr(smart_chunk, 'section_hierarchy') and smart_chunk.section_hierarchy:
3277:         hierarchy = smart_chunk.section_hierarchy
3278:         if len(hierarchy) > 0:
3279:             axes.append(hierarchy[0]) # Top-level = axis
3280:         if len(hierarchy) > 1:
3281:             programs.append(hierarchy[1]) # Second level = program
3282:         if len(hierarchy) > 2:
3283:             projects.append(hierarchy[2]) # Third level = project
3284:
3285:     return PolicyFacet(
3286:         axes=axes[:3], # Limit to avoid bloat
3287:         programs=programs[:5],
3288:         projects=projects[:5]
3289:     )
3290:
3291: @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_time_facets")
3292: def _extract_time_facets(self, smart_chunk: Any) -> TimeFacet:
3293:     """Extract temporal information from SPC temporal_dynamics."""
3294:     years = []
3295:     periods = []
3296:
3297:     if hasattr(smart_chunk, 'temporal_dynamics') and smart_chunk.temporal_dynamics:
3298:         temp = smart_chunk.temporal_dynamics
3299:         # Extract years from temporal_markers
3300:         if hasattr(temp, 'temporal_markers'):
3301:             for marker in temp.temporal_markers[:10]:
3302:                 # marker format: (text, marker_type, position)
3303:                 marker_text = marker[0] if isinstance(marker, (list, tuple)) else str(marker)
3304:                 # Try to extract years (4-digit numbers between 2020-2030)
3305:                 import re
3306:                 year_matches = re.findall(r'\b(202[0-9]|203[0-9])\b', marker_text)
3307:                 years.extend(int(y) for y in year_matches)
3308:
3309:     # Also check strategic_context for temporal_horizon
3310:     if hasattr(smart_chunk, 'strategic_context') and smart_chunk.strategic_context:
3311:         ctx = smart_chunk.strategic_context
3312:         if hasattr(ctx, 'temporal_horizon'):
3313:             periods.append(ctx.temporal_horizon)
3314:
3315:     return TimeFacet(
3316:         years=sorted(set(years))[:10], # Unique and sorted
3317:         periods=periods[:5]

```

```

3318:         )
3319:
3320:     @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_geo_facets")
3321:     def _extract_geo_facets(self, smart_chunk: Any) -> GeoFacet:
3322:         """Extract geographic information from SPC strategic_context."""
3323:         territories = []
3324:         regions = []
3325:
3326:         if hasattr(smart_chunk, 'strategic_context') and smart_chunk.strategic_context:
3327:             ctx = smart_chunk.strategic_context
3328:             if hasattr(ctx, 'geographic_scope'):
3329:                 territories.append(ctx.geographic_scope)
3330:                 # Could also parse from policy_entities with location types
3331:
3332:             return GeoFacet(
3333:                 territories=territories[:5],
3334:                 regions=regions[:5]
3335:             )
3336:
3337:     @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._build_provenance")
3338:     def _build_provenance(self, smart_chunk: Any) -> ProvenanceMap:
3339:         """Build provenance from SPC metadata."""
3340:         # Extract section info from section_hierarchy
3341:         source_section = None
3342:         if hasattr(smart_chunk, 'section_hierarchy') and smart_chunk.section_hierarchy:
3343:             source_section = " > ".join(smart_chunk.section_hierarchy[:3])
3344:
3345:         return ProvenanceMap(
3346:             source_page=None, # SPC doesn't track page numbers
3347:             source_section=source_section,
3348:             extraction_method="smart_policy_chunking_v3.0"
3349:         )
3350:
3351:     @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_entities")
3352:     def _extract_entities(self, smart_chunk: Any) -> list[Entity]:
3353:         """Extract entities from SPC policy_entities."""
3354:         entities = []
3355:
3356:         if hasattr(smart_chunk, 'policy_entities') and smart_chunk.policy_entities:
3357:             for pe in smart_chunk.policy_entities[:20]: # Limit to 20
3358:                 entity = Entity(
3359:                     text=pe.text if hasattr(pe, 'text') else str(pe),
3360:                     entity_type=pe.entity_type if hasattr(pe, 'entity_type') else 'unknown',
3361:                     confidence=pe.confidence if hasattr(pe, 'confidence') else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartCh
unkConverter._extract_entities", "auto_param_I373_79", 0.8)
3362:                 )
3363:                 entities.append(entity)
3364:
3365:         return entities
3366:
3367:     @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_budget")
3368:     def _extract_budget(self, smart_chunk: Any) -> Budget | None:
3369:         """
3370:             Extract budget with comprehensive error handling and logging (H1.4).
3371:
3372:             Implements 4 regex patterns and robust year extraction with 4 fallback strategies.

```

```
3373:     """
3374:     if not (hasattr(smart_chunk, 'strategic_context') and smart_chunk.strategic_context):
3375:         return None
3376:
3377:     ctx = smart_chunk.strategic_context
3378:     if not (hasattr(ctx, 'budget_linkage') and ctx.budget_linkage):
3379:         return None
3380:
3381:     import re
3382:     budget_text = ctx.budget_linkage
3383:
3384:     # H1.4: 4 regex patterns for robust budget extraction
3385:     patterns = [
3386:         # Pattern 1: Currency symbol with optional scale
3387:         r'\$\?s*([0-9,]+(?:\.[0-9]+)?)\s*(millones|mil millones|billion|billones)?',
3388:         # Pattern 2: "X millones de pesos"
3389:         r'(\d+(?:\.\d+)?)\s*millones?\s+de\s+pesos',
3390:         # Pattern 3: COP currency code
3391:         r'COP\s*\$?\s*([0-9,]+(?:\.[0-9]+)?)',
3392:         # Pattern 4: "presupuesto de $X"
3393:         r'presupuesto\s+de\s+\$?\s*([0-9,]+(?:\.[0-9]+)?)',
3394:     ]
3395:
3396:     amount = None
3397:     scale_multiplier = 1
3398:
3399:     for pattern in patterns:
3400:         match = re.search(pattern, budget_text, re.IGNORECASE)
3401:         if match:
3402:             try:
3403:                 amount_str = match.group(1).replace(',', '')
3404:                 amount = float(amount_str)
3405:
3406:                 # Detect scale from second group or text
3407:                 scale_text = budget_text.lower()
3408:                 if 'millones' in scale_text or 'million' in scale_text:
3409:                     scale_multiplier = 1_000_000
3410:                 elif 'mil millones' in scale_text or 'billion' in scale_text or 'billones' in scale_text:
3411:                     scale_multiplier = 1_000_000_000
3412:
3413:                 amount *= scale_multiplier
3414:                 break # Stop at first match
3415:
3416:             except (ValueError, IndexError) as e:
3417:                 self.logger.warning(f"Budget pattern matched but parsing failed: {e}")
3418:                 continue
3419:
3420:             if amount is None:
3421:                 return None
3422:
3423:     # H1.4: Extract year with 4 fallback strategies
3424:     year = self._extract_budget_year(budget_text, smart_chunk)
3425:
3426:     # Build Budget object
3427:     use = ctx.policy_intent if hasattr(ctx, 'policy_intent') else "General"
3428:
```

```
3429:         self.logger.debug(
3430:             f"Extracted budget: amount={amount:.2f}, year={year}, "
3431:             f"source='Strategic Context', use='{use[:30]}'"
3432:         )
3433:
3434:     return Budget(
3435:         source="Strategic Context",
3436:         use=use,
3437:         amount=amount,
3438:         year=year,
3439:         currency="COP"
3440:     )
3441:
3442: @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_budget_year")
3443: def _extract_budget_year(self, budget_text: str, smart_chunk: Any) -> int:
3444:     """
3445:     Extract budget year with 4 fallback strategies (H1.4).
3446:
3447:     Strategy 1: Extract from budget_text itself
3448:     Strategy 2: Extract from temporal_dynamics markers
3449:     Strategy 3: Use temporal_horizon from strategic_context
3450:     Strategy 4: Default to 2024
3451:     """
3452:     import re
3453:
3454:     # Strategy 1: Look for year in budget_text
3455:     year_match = re.search(r'\b(202[0-9]|203[0-9])\b', budget_text)
3456:     if year_match:
3457:         return int(year_match.group(1))
3458:
3459:     # Strategy 2: Check temporal_dynamics markers
3460:     if hasattr(smart_chunk, 'temporal_dynamics') and smart_chunk.temporal_dynamics:
3461:         temp = smart_chunk.temporal_dynamics
3462:         if hasattr(temp, 'temporal_markers') and temp.temporal_markers:
3463:             for marker in temp.temporal_markers:
3464:                 marker_text = marker[0] if isinstance(marker, (list, tuple)) else str(marker)
3465:                 year_match = re.search(r'\b(202[0-9]|203[0-9])\b', marker_text)
3466:                 if year_match:
3467:                     return int(year_match.group(1))
3468:
3469:     # Strategy 3: Check strategic_context temporal_horizon
3470:     if hasattr(smart_chunk, 'strategic_context') and smart_chunk.strategic_context:
3471:         ctx = smart_chunk.strategic_context
3472:         if hasattr(ctx, 'temporal_horizon') and ctx.temporal_horizon:
3473:             horizon_match = re.search(r'\b(202[0-9]|203[0-9])\b', ctx.temporal_horizon)
3474:             if horizon_match:
3475:                 return int(horizon_match.group(1))
3476:
3477:     # Strategy 4: Default to 2024
3478:     self.logger.debug("No year found in budget context, defaulting to 2024")
3479:     return 2024
3480:
3481: @calibrated_method("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi")
3482: def _extract_kpi(self, smart_chunk: Any) -> KPI | None:
3483:     """Extract KPI if chunk contains indicator information."""
3484:     # Check if chunk_type suggests this is a metric
```

```

3485:     chunk_type_str = smart_chunk.chunk_type.value if hasattr(smart_chunk.chunk_type, 'value') else str(smart_chunk.chunk_type)
3486:
3487:     if 'METRICA' in chunk_type_str.upper() or 'EVALUACION' in chunk_type_str.upper():
3488:         # Try to extract indicator from text
3489:         import re
3490:         text = smart_chunk.text
3491:         # Look for percentage targets (e.g., "90%", "meta del 85%")
3492:         target_match = re.search(r'meta.*?([0-9]+(?:\.[0-9]+)?)\s*%', text, re.IGNORECASE)
3493:         if target_match:
3494:             return KPI(
3495:                 indicator_name=smart_chunk.text[:80], # Use chunk text as indicator name
3496:                 target_value=float(target_match.group(1)),
3497:                 unit="%",
3498:                 year=None
3499:             )
3500:
3501:     return None
3502:
3503: def _calculate_quality_metrics(
3504:     self,
3505:     smart_chunks: list[Any],
3506:     chunk_graph: ChunkGraph
3507: ) -> QualityMetrics:
3508:     """Calculate quality metrics from SPC data."""
3509:     # Provenance completeness
3510:     chunks_with_provenance = sum(
3511:         1 for c in chunk_graph.chunks.values()
3512:             if c.provenance and c.provenance.source_section
3513:     )
3514:     provenance_completeness = chunks_with_provenance / len(chunk_graph.chunks) if chunk_graph.chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L526_110", 0.0)
3515:
3516:     # Average coherence from SPC coherence_score
3517:     avg_coherence = sum(sc.coherence_score for sc in smart_chunks) / len(smart_chunks) if smart_chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L529_112", 0.0)
3518:
3519:     # Average completeness from SPC completeness_index
3520:     avg_completeness = sum(sc.completeness_index for sc in smart_chunks) / len(smart_chunks) if smart_chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L532_118", 0.0)
3521:
3522:     # Budget consistency
3523:     chunks_with_budget = sum(1 for c in chunk_graph.chunks.values() if c.budget)
3524:     budget_consistency = chunks_with_budget / len(chunk_graph.chunks) if chunk_graph.chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L536_101", 0.0)
3525:
3526:     # Temporal robustness
3527:     chunks_with_time = sum(1 for c in chunk_graph.chunks.values() if c.time_facets.years)
3528:     temporal_robustness = chunks_with_time / len(chunk_graph.chunks) if chunk_graph.chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L540_100", 0.0)
3529:
3530:     # Chunk context coverage (from edges)
3531:     chunks_with_edges = len({e[0] for e in chunk_graph.edges} | {e[1] for e in chunk_graph.edges})
3532:     chunk_context_coverage = chunks_with_edges / len(chunk_graph.chunks) if chunk_graph.chunks else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L544_104", 0.0)
3533:
3534:     return QualityMetrics(

```

```

3535:         provenance_completeness=provenance_completeness,
3536:         structural_consistency=avg_coherence,
3537:         boundary_f1=avg_completeness,
3538:         kpi_linkage_rate=ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.SmartChunkConverter._extract_kpi", "auto_param_L550_29",
0.0), # Would need KPI analysis
3539:         budget_consistency_score=budget_consistency,
3540:         temporal_robustness=temporal_robustness,
3541:         chunk_context_coverage=chunk_context_coverage
3542:     )
3543:
3544:     def _generate_integrity_index(
3545:         self,
3546:         chunk_hashes: dict[str, str],
3547:         document_metadata: dict[str, Any]
3548:     ) -> IntegrityIndex:
3549:         """
3550:             Generate cryptographic integrity index.
3551:
3552:             Uses BLAKE2b-256 to compute aggregate hash of all chunk hashes.
3553:             NOT a true Merkle tree - simply hashes sorted JSON representation.
3554:         """
3555:         # Generate root hash from all chunk hashes (sorted for determinism)
3556:         combined = json.dumps(chunk_hashes, sort_keys=True).encode('utf-8')
3557:         blake2b_root = hashlib.blake2b(combined, digest_size=32).hexdigest()
3558:
3559:         return IntegrityIndex(
3560:             blake2b_root=blake2b_root,
3561:             chunk_hashes=chunk_hashes
3562:         )
3563:
3564:     def _preserve_spc_rich_data(
3565:         self,
3566:         smart_chunks: list[Any],
3567:         document_metadata: dict[str, Any]
3568:     ) -> dict[str, Any]:
3569:         """
3570:             Preserve SPC rich data in metadata for executor access.
3571:
3572:             This is critical for enabling executors to access the full SPC analysis:
3573:             - Embeddings (semantic, policy, causal, temporal)
3574:             - Causal chains with evidence
3575:             - Strategic context
3576:             - Quality scores
3577:         """
3578:         enriched = dict(document_metadata)
3579:
3580:         # Store serializable SPC data
3581:         spc_rich_data = {}
3582:
3583:         for sc in smart_chunks:
3584:             chunk_data = {
3585:                 'chunk_id': sc.chunk_id,
3586:                 'semantic_density': sc.semantic_density,
3587:                 'coherence_score': sc.coherence_score,
3588:                 'completeness_index': sc.completeness_index,
3589:                 'strategic_importance': sc.strategic_importance,

```

```

3590:             'information_density': sc.information_density,
3591:             'actionability_score': sc.actionability_score,
3592:         }
3593:
3594:     # Add embeddings if available (as lists for JSON serialization)
3595:     # CRITICAL: Fail-fast if embeddings cannot be preserved (no silent data loss)
3596:     if hasattr(sc, 'semantic_embedding') and sc.semantic_embedding is not None:
3597:         try:
3598:             import numpy as np
3599:         except ImportError as e:
3600:             self.logger.error(
3601:                 f"Chunk {sc.chunk_id}: NumPy is required for embedding preservation but not available"
3602:             )
3603:             raise RuntimeError(
3604:                 "NumPy is required for SPC embedding preservation. "
3605:                 "Install with: pip install numpy>=1.26.0"
3606:             ) from e
3607:
3608:     # Validate embedding type
3609:     if not isinstance(sc.semantic_embedding, np.ndarray):
3610:         self.logger.error(
3611:             f"Chunk {sc.chunk_id}: semantic_embedding is not np.ndarray, "
3612:             f"got {type(sc.semantic_embedding)}"
3613:         )
3614:         raise TypeError(
3615:             f"Expected semantic_embedding to be np.ndarray, got {type(sc.semantic_embedding)}"
3616:         )
3617:
3618:     # Convert to list for JSON serialization
3619:     try:
3620:         chunk_data['semantic_embedding'] = sc.semantic_embedding.tolist()
3621:         chunk_data['embedding_dim'] = sc.semantic_embedding.shape[0]
3622:         self.logger.debug(
3623:             f"Chunk {sc.chunk_id}: Preserved embedding with dimension {sc.semantic_embedding.shape[0]}"
3624:         )
3625:     except (AttributeError, IndexError) as e:
3626:         self.logger.error(
3627:             f"Chunk {sc.chunk_id}: Failed to convert embedding to list: {e}"
3628:         )
3629:         raise RuntimeError(
3630:             f"Embedding conversion failed for chunk {sc.chunk_id}: {e}"
3631:         ) from e
3632:
3633:     # Add causal chain summary
3634:     if hasattr(sc, 'causal_chain') and sc.causal_chain:
3635:         chunk_data['causal_chain_count'] = len(sc.causal_chain)
3636:         chunk_data['causal_evidence'] = [
3637:             {
3638:                 'dimension': ce.dimension if hasattr(ce, 'dimension') else 'unknown',
3639:                 'confidence': ce.confidence if hasattr(ce, 'confidence') else ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.converter.
SmartChunkConverter._extract_kpi", "auto_param_L651_86", 0.0),
3640:             }
3641:             for ce in sc.causal_chain[:5]  # Top 5
3642:         ]
3643:
3644:     # Add strategic context summary

```

```
3645:         if hasattr(sc, 'strategic_context') and sc.strategic_context:
3646:             ctx = sc.strategic_context
3647:             chunk_data['strategic_context'] = {
3648:                 'policy_intent': ctx.policy_intent if hasattr(ctx, 'policy_intent') else None,
3649:                 'implementation_phase': ctx.implementation_phase if hasattr(ctx, 'implementation_phase') else None,
3650:                 'temporal_horizon': ctx.temporal_horizon if hasattr(ctx, 'temporal_horizon') else None,
3651:             }
3652:
3653:             # Add topic distribution
3654:             if hasattr(sc, 'topic_distribution') and sc.topic_distribution:
3655:                 chunk_data['topic_distribution'] = dict(sc.topic_distribution)
3656:
3657:             spc_rich_data[sc.chunk_id] = chunk_data
3658:
3659:             enriched['spc_rich_data'] = spc_rich_data
3660:             enriched['spc_version'] = 'SMART-CHUNK-3.0-FINAL'
3661:             enriched['conversion_timestamp'] = document_metadata.get('processing_timestamp', 'unknown')
3662:
3663:             self.logger.info(f"Preserved rich SPC data for {len(spc_rich_data)} chunks")
3664:
3665:             return enriched
3666:
3667:
3668:
3669: =====
3670: FILE: src/farfan_pipeline/processing/embedding_policy.py
3671: =====
3672:
3673: """
3674: INTERNAL SPC COMPONENT
3675:
3676: \232 \217 USAGE RESTRICTION \232 \217
3677: =====
3678: This module implements SOTA semantic embedding and policy analysis for Smart
3679: Policy Chunks. It MUST NOT be used as a standalone ingestion pipeline in the
3680: canonical FARFAN flow.
3681:
3682: Canonical entrypoint is scripts/run_policy_pipeline_verified.py.
3683:
3684: This module is an INTERNAL COMPONENT of:
3685:     src/farfan_core/processing/spc_ingestion.py (StrategicChunkingSystem)
3686:
3687: DO NOT use this module directly as an independent pipeline. It is consumed
3688: internally by the SPC core and should only be imported from within:
3689:     - farfan_core.processing.spc_ingestion
3690:     - Unit tests for SPC components
3691:
3692: State-of-the-Art Components:
3693:     - BGE-M3 multilingual embeddings (2024 SOTA)
3694:     - Cross-encoder reranking for Spanish policy documents
3695:     - Bayesian uncertainty quantification for numerical analysis
3696:     - Graph-based multi-hop reasoning
3697: =====
3698: """
3699:
3700: from __future__ import annotations
```

```
3701:  
3702: import hashlib  
3703: import logging  
3704: import re  
3705: from dataclasses import dataclass  
3706: from enum import Enum  
3707: from functools import lru_cache  
3708: from typing import TYPE_CHECKING, Any, Literal, Protocol, TypedDict  
3709:  
3710: import numpy as np  
3711: from sentence_transformers import CrossEncoder, SentenceTransformer  
3712: from sklearn.metrics.pairwise import cosine_similarity  
3713: from farfan_pipeline.core.parameters import ParameterLoaderV2  
3714: from farfan_pipeline.core.calibration.decorators import calibrated_method  
3715:  
3716: if TYPE_CHECKING:  
3717:     from collections.abc import Iterable  
3718:  
3719:     from numpy.typing import NDArray  
3720:  
3721: # =====  
3722: # DESIGN CONSTANTS - Model Configuration  
3723: # =====  
3724:  
3725: # Model constants  
3726: DEFAULT_CROSS_ENCODER_MODEL = "cross-encoder/ms-marco-MiniLM-L-6-v2"  
3727: MODEL_PARAPHRASE_MULTILINGUAL = "sentence-transformers/paraphrase-multilingual-mpnet-base-v2"  
3728:  
3729: # =====  
3730: # TYPE SYSTEM - Python 3.10+ Type Safety  
3731: # =====  
3732:  
3733: class PolicyDomain(Enum):  
3734:     """  
3735:         Colombian PDM policy areas (PA01-PA10) per canonical notation.  
3736:  
3737:         Values are loaded from questionnaire_monolith.json canonical_notation.  
3738:         Use CanonicalPolicyArea from farfan_core.core.canonical_notation for dynamic access.  
3739:         """  
3740:  
3741:         # Legacy IDs mapped to canonical codes for backward compatibility  
3742:         P1 = "PA01" # Derechos de las mujeres e igualdad de gÃ©nero  
3743:         P2 = "PA02" # PrevenciÃ³n de la violencia y protecciÃ³n frente al conflicto  
3744:         P3 = "PA03" # Ambiente sano, cambio climÃ¡tico, prevenciÃ³n y atenciÃ³n a desastres  
3745:         P4 = "PA04" # Derechos econÃ³micos, sociales y culturales  
3746:         P5 = "PA05" # Derechos de las vÃ—ctimas y construcciÃ³n de paz  
3747:         P6 = "PA06" # Derecho al buen futuro de la niÃ±ez, adolescencia, juventud  
3748:         P7 = "PA07" # Tierras y territorios  
3749:         P8 = "PA08" # LÃ—deres y defensores de derechos humanos  
3750:         P9 = "PA09" # Crisis de derechos de personas privadas de la libertad  
3751:         P10 = "PA10" # MigraciÃ³n transfronteriza  
3752:  
3753: class AnalyticalDimension(Enum):  
3754:     """  
3755:         Analytical dimensions (D1-D6) per canonical notation.  
3756:
```

```
3757:     Values reference canonical notation from questionnaire_monolith.json.
3758:     Use CanonicalDimension from farfan_core.core.canonical_notation for dynamic access.
3759:     """
3760:
3761:     D1 = "DIM01" # INSUMOS - DiagnÃstico y Recursos
3762:     D2 = "DIM02" # ACTIVIDADES - DiseÃn de IntervenciÃn
3763:     D3 = "DIM03" # PRODUCTOS - Productos y Outputs
3764:     D4 = "DIM04" # RESULTADOS - Resultados y Outcomes
3765:     D5 = "DIM05" # IMPACTOS - Impactos de Largo Plazo
3766:     D6 = "DIM06" # CAUSALIDAD - TeorÃa de Cambio
3767:
3768: class PDQIdentifier(TypedDict):
3769:     """Canonical P-D-Q identifier structure."""
3770:
3771:     question_unique_id: str # P#-D#-Q#
3772:     policy: str # P#
3773:     dimension: str # D#
3774:     question: int # Q#
3775:     rubric_key: str # D#-Q#
3776:
3777: class PosteriorSampleRecord(TypedDict):
3778:     """Serializable posterior sample used by downstream Bayesian consumers."""
3779:
3780:     coherence: float
3781:
3782: class SemanticChunk(TypedDict):
3783:     """Structured semantic chunk with metadata."""
3784:
3785:     chunk_id: str
3786:     content: str
3787:     embedding: NDArray[np.float32]
3788:     metadata: dict[str, Any]
3789:     pdq_context: PDQIdentifier | None
3790:     token_count: int
3791:     position: tuple[int, int] # (start, end) in document
3792:
3793: class PosteriorSample(TypedDict):
3794:     """Serialized posterior sample representation."""
3795:
3796:     coherence: float
3797:
3798: class BayesianEvaluation(TypedDict):
3799:     """Bayesian uncertainty-aware evaluation result."""
3800:
3801:     point_estimate: float # 0.0-1.0
3802:     credible_interval_95: tuple[float, float]
3803:     posterior_samples: list[PosteriorSample]
3804:     evidence_strength: Literal["weak", "moderate", "strong", "very_strong"]
3805:     numerical_coherence: float # Statistical consistency score
3806:     posterior_records: list[PosteriorSampleRecord]
3807:
3808: class EmbeddingProtocol(Protocol):
3809:     """Protocol for embedding models."""
3810:
3811:     def encode(
3812:         self, texts: list[str], batch_size: int = 32, normalize: bool = True
```

```
3813:     ) -> NDArray[np.float32]: ...
3814:
3815: def to_dict_samples(samples: NDArray[np.float32] | Iterable[float]) -> list[PosteriorSample]:
3816:     """Convert posterior samples to the serialized TypedDict format."""
3817:
3818:     array = np.asarray(list(samples)) if not hasattr(samples, "shape") else samples, dtype=np.float32)
3819:     flat = array.ravel()
3820:     return [{"coherence": float(value)} for value in flat]
3821:
3822: def samples_to_array(samples: NDArray[np.float32] | Iterable[PosteriorSample]) -> NDArray[np.float32]:
3823:     """Normalize posterior samples into a numpy array for computation."""
3824:
3825:     if isinstance(samples, np.ndarray):
3826:         return samples.astype(np.float32)
3827:     return np.array([sample["coherence"] for sample in samples], dtype=np.float32)
3828:
3829: def ensure_content_schema(chunk: dict[str, Any]) -> dict[str, Any]:
3830:     """Ensure chunk dictionaries expose the ``content`` key."""
3831:
3832:     if "content" not in chunk and "text" in chunk:
3833:         upgraded = dict(chunk)
3834:         upgraded["content"] = upgraded.pop("text")
3835:         return upgraded
3836:     return chunk
3837:
3838: # =====
3839: # ADVANCED SEMANTIC CHUNKING - State-of-the-Art
3840: # =====
3841:
3842: @dataclass
3843: class ChunkingConfig:
3844:     """Configuration for semantic chunking optimized for PDM documents."""
3845:
3846:     chunk_size: int = 512 # Tokens, optimized for policy documents
3847:     chunk_overlap: int = 128 # Preserve context across chunks
3848:     min_chunk_size: int = 64 # Avoid tiny fragments
3849:     respect_boundaries: bool = True # Sentence/paragraph boundaries
3850:     preserve_tables: bool = True # Keep tables intact
3851:     detect_lists: bool = True # Recognize enumerations
3852:     section_aware: bool = True # Understand document structure
3853:
3854: class AdvancedSemanticChunker:
3855:     """
3856:         State-of-the-art semantic chunking for Colombian policy documents.
3857:
3858:     Implements:
3859:         - Recursive character splitting with semantic boundary preservation
3860:         - Table structure detection and preservation
3861:         - List and enumeration recognition
3862:         - Hierarchical section awareness (P-D-Q structure)
3863:         - Token-aware splitting (not just character-based)
3864:     """
3865:
3866:     # Colombian policy document patterns
3867:     SECTION_HEADERS = re.compile(
3868:         r"^(?:CAPÃ\215TULO|SECCIÃ\223N|ARTÃ\215CULO|PROGRAMA|PROYECTO|EJE)\s+[IVX\d]+",
```

```
3869:         re.MULTILINE | re.IGNORECASE,
3870:     )
3871:     TABLE_MARKERS = re.compile(r"(?:Tabla|Cuadro|Figura)\s+\d+", re.IGNORECASE)
3872:     LIST_MARKERS = re.compile(r"\s*\n[\s]*[\u200c-\u202f]\d+[\.\.])\s+", re.MULTILINE)
3873:     NUMERIC_INDICATORS = re.compile(
3874:         r"\b\d+(?:[.,]\d+)?(?:\s*%|millones?|mil|billones?)?\b", re.IGNORECASE
3875:     )
3876:
3877:     def __init__(self, config: ChunkingConfig) -> None:
3878:         self.config = config
3879:         self._logger = logging.getLogger(self.__class__.__name__)
3880:
3881:     def chunk_document(
3882:         self,
3883:         *,
3884:         text: str,
3885:         document_metadata: dict[str, Any],
3886:     ) -> list[SemanticChunk]:
3887:         """
3888:             Chunk document with advanced semantic awareness (keyword-only params).
3889:
3890:             Args:
3891:                 text: Document text to chunk
3892:                 document_metadata: Metadata dict with at least 'doc_id' key
3893:
3894:             Returns:
3895:                 List of semantic chunks with preserved structure and P-D-Q context
3896:
3897:             Raises:
3898:                 TypeError: If text is not a string
3899:                 KeyError: If document_metadata missing required keys
3900:         """
3901:         # Runtime validation at ingress
3902:         if not isinstance(text, str):
3903:             raise TypeError(
3904:                 f"ERR_CONTRACT_MISMATCH[fn=chunk_document, param='text', "
3905:                 f"expected=str, got={type(text).__name__}]"
3906:             )
3907:
3908:         if not isinstance(document_metadata, dict):
3909:             raise TypeError(
3910:                 f"ERR_CONTRACT_MISMATCH[fn=chunk_document, param='document_metadata', "
3911:                 f"expected=dict, got={type(document_metadata).__name__}]"
3912:             )
3913:         # Preprocess: normalize whitespace, preserve structure
3914:         normalized_text = self._normalize_text(text)
3915:
3916:         # Extract structural elements
3917:         sections = self._extract_sections(normalized_text)
3918:         tables = self._extract_tables(normalized_text)
3919:         lists = self._extract_lists(normalized_text)
3920:
3921:         # Generate chunks with boundary preservation
3922:         raw_chunks = self._recursive_split(
3923:             normalized_text,
3924:             target_size=self.config.chunk_size,
```

```

3925:         overlap=self.config.chunk_overlap,
3926:     )
3927:
3928:     # Enrich chunks with metadata and P-D-Q context
3929:     semantic_chunks: list[SemanticChunk] = []
3930:
3931:     for idx, chunk_text in enumerate(raw_chunks):
3932:         # Infer P-D-Q context from chunk text
3933:         pdq_context = self._infer_pdq_context(chunk_text)
3934:
3935:         # Count tokens (approximation: Spanish has ~1.3 chars/token)
3936:         AVG_CHARS_PER_TOKEN = 1.3 # Source: Spanish language statistics
3937:         token_count = int(
3938:             len(chunk_text) / AVG_CHARS_PER_TOKEN
3939:         ) # Approximate token count
3940:
3941:         # Create structured chunk
3942:         chunk_id = hashlib.sha256(
3943:             f"{{document_metadata.get('doc_id', '')}}_{{idx}}_{{chunk_text[:50]}}".encode()
3944:         ).hexdigest()[:16]
3945:
3946:         semantic_chunk: SemanticChunk = {
3947:             "chunk_id": chunk_id,
3948:             "content": chunk_text,
3949:             "embedding": np.array([]), # Filled later
3950:             "metadata": (
3951:                 "document_id": document_metadata.get("doc_id"),
3952:                 "chunk_index": idx,
3953:                 "has_table": self._contains_table(chunk_text, tables),
3954:                 "has_list": self._contains_list(chunk_text, lists),
3955:                 "has_numbers": bool(self.NUMERIC_INDICATORS.search(chunk_text)),
3956:                 "section_title": self._find_section(chunk_text, sections),
3957:             ),
3958:             "pdq_context": pdq_context,
3959:             "token_count": token_count,
3960:             "position": (0, len(chunk_text)), # Updated during splitting
3961:         }
3962:
3963:         semantic_chunks.append(ensure_content_schema(semantic_chunk))
3964:
3965:         self._logger.info(
3966:             "Created %d semantic chunks from document %s",
3967:             len(semantic_chunks),
3968:             document_metadata.get("doc_id", "unknown"),
3969:         )
3970:
3971:         return semantic_chunks
3972:
3973: @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._normalize_text")
3974: def _normalize_text(self, text: str) -> str:
3975:     """Normalize text while preserving structure."""
3976:     # Remove excessive whitespace but preserve paragraph breaks
3977:     text = re.sub(r"\t+", " ", text)
3978:     text = re.sub(r"\n{3,}", "\n\n", text)
3979:     return text.strip()
3980:

```

```

3981:     @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._recursive_split")
3982:     def _recursive_split(self, text: str, target_size: int, overlap: int) -> list[str]:
3983:         """
3984:             Recursive character splitting with semantic boundary respect.
3985:
3986:             Priority: Paragraph > Sentence > Word > Character
3987:             """
3988:             if len(text) <= target_size:
3989:                 return [text]
3990:
3991:             chunks = []
3992:             current_pos = 0
3993:
3994:             while current_pos < len(text):
3995:                 # Calculate chunk end position
3996:                 end_pos = min(current_pos + target_size, len(text))
3997:
3998:                 # Try to find semantic boundary
3999:                 if end_pos < len(text):
4000:                     # Priority 1: Paragraph break
4001:                     paragraph_break = text.rfind("\n\n", current_pos, end_pos)
4002:                     if paragraph_break != -1 and paragraph_break > current_pos:
4003:                         end_pos = paragraph_break + 2
4004:
4005:                     # Priority 2: Sentence boundary
4006:                     elif sentence_end := self._find_sentence_boundary(
4007:                         text, current_pos, end_pos
4008:                     ):
4009:                         end_pos = sentence_end
4010:
4011:                     chunk = text[current_pos:end_pos].strip()
4012:                     if len(chunk) >= self.config.min_chunk_size:
4013:                         chunks.append(chunk)
4014:
4015:                     # Move position with overlap
4016:                     current_pos = end_pos - overlap if overlap > 0 else end_pos
4017:
4018:                     # Prevent infinite loop
4019:                     if current_pos <= end_pos - target_size:
4020:                         current_pos = end_pos
4021:
4022:             return chunks
4023:
4024:     @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._find_sentence_boundary")
4025:     def _find_sentence_boundary(self, text: str, start: int, end: int) -> int | None:
4026:         """Find sentence boundary using Spanish punctuation rules."""
4027:         # Spanish sentence endings: . ! ? ; followed by space or newline
4028:         sentence_pattern = re.compile(r"[.!?;]\s+")
4029:
4030:         matches = list(sentence_pattern.finditer(text, start, end))
4031:         if matches:
4032:             # Return position after punctuation and space
4033:             return matches[-1].end()
4034:         return None
4035:
4036:     @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._extract_sections")

```

```
4037:     def _extract_sections(self, text: str) -> list[dict[str, Any]]:
4038:         """Extract document sections with hierarchical structure."""
4039:         sections = []
4040:         for match in self.SECTION_HEADERS.finditer(text):
4041:             sections.append(
4042:                 {
4043:                     "title": match.group(0),
4044:                     "position": match.start(),
4045:                     "end": match.end(),
4046:                 }
4047:             )
4048:         return sections
4049:
4050: # Number of characters to consider as table extent after marker
4051: TABLE_EXTENT_CHARS = 300
4052:
4053: @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._extract_tables")
4054: def _extract_tables(self, text: str) -> list[dict[str, Any]]:
4055:     """Identify table regions in document."""
4056:     tables = []
4057:     for match in self.TABLE_MARKERS.finditer(text):
4058:         # Heuristic: table extends ~TABLE_EXTENT_CHARS chars after marker
4059:         tables.append(
4060:             {
4061:                 "marker": match.group(0),
4062:                 "start": match.start(),
4063:                 "end": min(match.end() + self.TABLE_EXTENT_CHARS, len(text)),
4064:             }
4065:         )
4066:     return tables
4067:
4068: @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._extract_lists")
4069: def _extract_lists(self, text: str) -> list[dict[str, Any]]:
4070:     """Identify list structures."""
4071:     lists = []
4072:     for match in self.LIST_MARKERS.finditer(text):
4073:         lists.append({"marker": match.group(0), "position": match.start()})
4074:     return lists
4075:
4076: def _infer_pdq_context(
4077:     self,
4078:     chunk_text: str,
4079: ) -> PDQIdentifier | None:
4080:     """
4081:     Infer P-D-Q context from chunk content.
4082:
4083:     Uses heuristics based on Colombian policy vocabulary.
4084:     """
4085:     # Policy-specific keywords (simplified for example)
4086:     policy_keywords = {
4087:         "PA01": ["mujer", "gÃ©nero", "igualdad", "equidad"],
4088:         "PA02": ["violencia", "conflicto", "seguridad", "prevenciÃ³n"],
4089:         "PA03": ["ambiente", "clima", "desastre", "riesgo"],
4090:         "PA04": ["econÃ³mico", "social", "cultural", "empleo"],
4091:         "PA05": ["vÃ-ctima", "paz", "reconciliaciÃ³n", "reparaciÃ³n"],
4092:         "PA06": ["niÃ±ez", "adolescente", "juventud", "futuro"],
```

```

4093:         "PA07": ["tierra", "territorio", "rural", "agrario"],
4094:         "PA08": ["lÃ-der", "defensor", "derechos humanos"],
4095:         "PA09": ["privado libertad", "cÃ;rcel", "reclusiÃ;n"],
4096:         "PA10": ["migraciÃ;n", "frontera", "venezolano"],
4097:     }
4098:
4099:     dimension_keywords = {
4100:         "DIM01": ["diagnÃ;stico", "baseline", "situaciÃ;n", "recurso"],
4101:         "DIM02": ["diseÃ;o", "estrategia", "intervenciÃ;n", "actividad"],
4102:         "DIM03": ["producto", "output", "entregable", "meta"],
4103:         "DIM04": ["resultado", "outcome", "efecto", "cambio"],
4104:         "DIM05": ["impacto", "largo plazo", "sostenibilidad"],
4105:         "DIM06": ["teorÃ;a", "causal", "coherencia", "lÃ;gica"],
4106:     }
4107:
4108:     # Score policies and dimensions
4109:     policy_scores = {
4110:         policy: sum(1 for kw in keywords if kw.lower() in chunk_text.lower())
4111:             for policy, keywords in policy_keywords.items()
4112:     }
4113:
4114:     dimension_scores = {
4115:         dim: sum(1 for kw in keywords if kw.lower() in chunk_text.lower())
4116:             for dim, keywords in dimension_keywords.items()
4117:     }
4118:
4119:     # Select best match if confidence is sufficient
4120:     best_policy = max(policy_scores, key=policy_scores.get)
4121:     best_dimension = max(dimension_scores, key=dimension_scores.get)
4122:
4123:     if policy_scores[best_policy] > 0 and dimension_scores[best_dimension] > 0:
4124:         # Generate canonical identifier
4125:         question_num = 1 # Simplified; real system would infer from context
4126:         question_code = f"Q{question_num:03d}"
4127:
4128:         return PDQIdentifier(
4129:             question_unique_id=f"{best_policy}-{best_dimension}-{question_code}",
4130:             policy=best_policy,
4131:             dimension=best_dimension,
4132:             question=question_num,
4133:             rubric_key=f"{best_dimension}-{question_code}",
4134:         )
4135:
4136:     return None
4137:
4138:     def _contains_table(
4139:         self, chunk_text: str, tables: list[dict[str, Any]]
4140:     ) -> bool:
4141:         """Check if chunk contains table markers."""
4142:         return any(
4143:             table["marker"] in chunk_text
4144:                 for table in tables
4145:         )
4146:
4147:     @calibrated_method("farfan_core.processing.embedding_policy.AdvancedSemanticChunker._contains_list")
4148:     def _contains_list(self, chunk_text: str, lists: list[dict[str, Any]]) -> bool:

```

```
4149:     """Check if chunk contains list structures."""
4150:     return bool(self.LIST_MARKERS.search(chunk_text))
4151:
4152:     def _find_section(
4153:         self, chunk_text: str, sections: list[dict[str, Any]]
4154:     ) -> str | None:
4155:         """Find section title for chunk."""
4156:         # Simplified: would use position-based matching in production
4157:         for section in sections:
4158:             if section["title"][:20] in chunk_text:
4159:                 return section["title"]
4160:
4161:         return None
4162: # =====
4163: # BAYESIAN NUMERICAL ANALYSIS - Rigorous Statistical Framework
4164: # =====
4165:
4166: class BayesianNumericalAnalyzer:
4167:     """
4168:         Bayesian framework for uncertainty-aware numerical policy analysis.
4169:
4170:         Implements:
4171:             - Beta-Binomial conjugate prior for proportions
4172:             - Normal-Normal conjugate prior for continuous metrics
4173:             - Bayesian hypothesis testing for policy comparisons
4174:             - Credible interval estimation
4175:             - Evidence strength quantification (Bayes factors)
4176:     """
4177:
4178:     def __init__(self, prior_strength: float = 1.0) -> None:
4179:         """
4180:             Initialize Bayesian analyzer.
4181:
4182:             Args:
4183:                 prior_strength: Prior belief strength (ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L510_51", 1.0) = weak, 1ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L510_64", 0.0) = strong)
4184:             """
4185:             self.prior_strength = prior_strength
4186:             self._logger = logging.getLogger(self.__class__.__name__)
4187:             self._rng = np.random.default_rng()
4188:
4189:             def evaluate_policy_metric(
4190:                 self,
4191:                 observed_values: list[float],
4192:                 n_posterior_samples: int = 10000,
4193:                 **kwargs: Any
4194:             ) -> BayesianEvaluation:
4195:                 """
4196:                     Bayesian evaluation of policy metric with uncertainty quantification.
4197:
4198:                     Returns posterior distribution, credible intervals, and evidence strength.
4199:
4200:                     Args:
4201:                         observed_values: List of observed metric values
4202:                         n_posterior_samples: Number of posterior samples to generate
```

```

4203:             **kwargs: Additional optional parameters for compatibility
4204:
4205:     Returns:
4206:         BayesianEvaluation with posterior samples and credible intervals
4207:
4208:     """
4209:     if not observed_values:
4210:         return self._null_evaluation()
4211:
4212:     obs_array = np.array(observed_values)
4213:
4214:     # Choose likelihood model based on data characteristics
4215:     if all(0 <= v <= 1 for v in observed_values):
4216:         # Proportion/probability metric: use Beta-Binomial
4217:         posterior_samples = self._beta_binomial_posterior(
4218:             obs_array, n_posterior_samples
4219:         )
4220:     else:
4221:         # Continuous metric: use Normal-Normal
4222:         posterior_samples = self._normal_normal_posterior(
4223:             obs_array, n_posterior_samples
4224:         )
4225:
4226:     # Compute statistics
4227:     point_estimate = float(np.median(posterior_samples))
4228:     ci_lower, ci_upper = (
4229:         float(np.percentile(posterior_samples, 2.5)),
4230:         float(np.percentile(posterior_samples, 97.5)),
4231:     )
4232:
4233:     # Quantify evidence strength using posterior width
4234:     ci_width = ci_upper - ci_lower
4235:     evidence_strength = self._classify_evidence_strength(ci_width)
4236:
4237:     # Assess numerical coherence (consistency of observations)
4238:     coherence = self._compute_coherence(obs_array)
4239:
4240:     serialized_samples = to_dict_samples(posterior_samples)
4241:
4242:     return BayesianEvaluation(
4243:         point_estimate=point_estimate,
4244:         credible_interval_95=(ci_lower, ci_upper),
4245:         posterior_samples=serialized_samples,
4246:         evidence_strength=evidence_strength,
4247:         numerical_coherence=coherence,
4248:         posterior_records=self.serialize_posterior_samples(posterior_samples),
4249:     )
4250:
4251:     def _beta_binomial_posterior(
4252:         self, observations: NDArray[np.float32], n_samples: int
4253:     ) -> NDArray[np.float32]:
4254:
4255:         """
4256:         Beta-Binomial conjugate posterior for proportion metrics.
4257:
4258:         Prior: Beta( $\hat{\pi}$ ,  $\hat{\pi}^2$ )
4259:         Likelihood: Binomial
4260:         Posterior: Beta( $\hat{\pi}$  + successes,  $\hat{\pi}^2$  + failures)

```

```

4259:     """
4260:     # Prior parameters (weakly informative)
4261:     alpha_prior = self.prior_strength
4262:     beta_prior = self.prior_strength
4263:
4264:     # Convert proportions to successes/failures
4265:     n_obs = len(observations)
4266:     sum_success = np.sum(observations)  # If already in [0,1]
4267:
4268:     # Posterior parameters
4269:     alpha_post = alpha_prior + sum_success
4270:     beta_post = beta_prior + (n_obs - sum_success)
4271:
4272:     # Sample from posterior
4273:     posterior_samples = self._rng.beta(alpha_post, beta_post, size=n_samples)
4274:
4275:     return posterior_samples.astype(np.float32)
4276:
4277: def _normal_normal_posterior(
4278:     self, observations: NDArray[np.float32], n_samples: int
4279: ) -> NDArray[np.float32]:
4280:     """
4281:     Normal-Normal conjugate posterior for continuous metrics.
4282:
4283:     Prior: Normal(1/202\200, 1/203\202\200\2)
4284:     Likelihood: Normal(1/4, 1/203\2)
4285:     Posterior: Normal(1/4_post, 1/203_post\2)
4286:     """
4287:     n_obs = len(observations)
4288:     obs_mean = np.mean(observations)
4289:     obs_std = np.std(observations, ddof=1) if n_obs > 1 else ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L616_65", 1.0)
4290:
4291:     # Prior parameters (weakly informative centered on observed mean)
4292:     mu_prior = obs_mean
4293:     sigma_prior = obs_std * self.prior_strength
4294:
4295:     # Posterior parameters (conjugate update)
4296:     precision_prior = 1 / (sigma_prior**2)
4297:     precision_likelihood = n_obs / (obs_std**2)
4298:
4299:     precision_post = precision_prior + precision_likelihood
4300:     mu_post = (
4301:         precision_prior * mu_prior + precision_likelihood * obs_mean
4302:     ) / precision_post
4303:     sigma_post = np.sqrt(1 / precision_post)
4304:
4305:     # Sample from posterior
4306:     posterior_samples = self._rng.normal(mu_post, sigma_post, size=n_samples)
4307:
4308:     return posterior_samples.astype(np.float32)
4309:
4310: def _classify_evidence_strength(
4311:     self, credible_interval_width: float, **kwargs: Any
4312: ) -> Literal["weak", "moderate", "strong", "very_strong"]:
4313:     """Classify evidence strength based on posterior uncertainty.

```

```
4314:
4315:     Args:
4316:         credible_interval_width: Width of the 95% credible interval
4317:         **kwargs: Additional optional parameters for compatibility
4318:
4319:     Returns:
4320:         Evidence strength classification (weak/moderate/strong/very_strong)
4321:         """
4322:         if credible_interval_width > ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L649_37",
4323:             0.5):
4324:             return "weak"
4325:         elif credible_interval_width > ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L651_39",
4326:             0.3):
4327:             return "moderate"
4328:         elif credible_interval_width > ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer.__init__", "auto_param_L653_39",
4329:             0.15):
4330:             return "strong"
4331:         else:
4332:             return "very_strong"
4333:
4334:     @calibrated_method("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._compute_coherence")
4335:     def _compute_coherence(self, observations: NDArray[np.float32], **kwargs: Any) -> float:
4336:
4337:         Compute numerical coherence (consistency) score.
4338:
4339:         Args:
4340:             observations: Array of observed values
4341:             **kwargs: Additional optional parameters for compatibility
4342:
4343:         Returns:
4344:             Coherence score in [0, 1]
4345:
4346:         if len(observations) < 2:
4347:             return ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._compute_coherence", "auto_param_L673_19", 1.0)
4348:
4349:             # Coefficient of variation
4350:             mean_val = np.mean(observations)
4351:             std_val = np.std(observations, ddof=1)
4352:
4353:             if mean_val == 0:
4354:                 return ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._compute_coherence", "auto_param_L680_19", 0.0)
4355:
4356:             cv = std_val / abs(mean_val)
4357:
4358:             # Normalize: lower CV = higher coherence
4359:             coherence = np.exp(-cv)  # Exponential decay
4360:
4361:             return float(np.clip(coherence, ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._compute_coherence", "auto_param_L687_40", 0.0), ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._compute_coherence", "auto_param_L687_45", 1.0)))
4362:
4363:     @calibrated_method("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation")
4364:     def _null_evaluation(self) -> BayesianEvaluation:
4365:         """Return null evaluation when no data available."""
4366:         null_samples = to_dict_samples(np.array([ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation",
```

```
"auto_param_L692_49", 0.0)], dtype=np.float32))
4366:
4367:     return BayesianEvaluation(
4368:         point_estimate=ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L695_27",
0.0),
4369:         credible_interval_95=(ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L696_34",
96_34), ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L696_39", 0.0)),
4370:         posterior_samples=null_samples,
4371:         evidence_strength="weak",
4372:         numerical_coherence=ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L699_32",
32), 0.0),
4373:         posterior_records=[{"coherence": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "au
to_param_L700_45", 0.0)}],
4374:     )
4375:
4376:     def serialize_posterior_samples(
4377:         self, samples: NDArray[np.float32]
4378:     ) -> list[PosteriorSampleRecord]:
4379:         """Convert posterior samples into standardized coherence records.
4380:
4381:             Safely handles None or non-array inputs and limits the number of
4382:             serialized records to avoid excessive memory use.
4383:         """
4384:         if samples is None:
4385:             return []
4386:
4387:         # Ensure a 1-D numpy array of floats
4388:         arr = np.asarray(samples, dtype=np.float32).ravel()
4389:
4390:         # Prevent accidental excessive memory use when serializing huge arrays
4391:         MAX_RECORDS = 10000
4392:         values = arr.tolist()
4393:         if len(values) > MAX_RECORDS:
4394:             values = values[:MAX_RECORDS]
4395:
4396:         return [{"coherence": float(v)} for v in values]
4397:
4398:     def compare_policies(
4399:         self,
4400:         policy_a_values: list[float],
4401:         policy_b_values: list[float],
4402:     ) -> dict[str, Any]:
4403:         """
4404:             Bayesian comparison of two policy metrics.
4405:
4406:             Returns probability that A > B and Bayes factor.
4407:         """
4408:         if not policy_a_values or not policy_b_values:
4409:             return {"probability_a_better": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "aut
o_param_L736_44", 0.5), "bayes_factor": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L73
6_65", 1.0)}
4410:
4411:         # Get posterior distributions
4412:         eval_a = self.evaluate_policy_metric(policy_a_values)
4413:         eval_b = self.evaluate_policy_metric(policy_b_values)
4414:
```

```

4415:     # Compute probability that A > B and clip to avoid exact 0/1 which can cause
4416:     # division-by-zero in subsequent Bayes factor calculation
4417:     samples_a = samples_to_array(eval_a["posterior_samples"])
4418:     samples_b = samples_to_array(eval_b["posterior_samples"])
4419:
4420:     # Compute probability that A > B and clip to avoid exact 0/1 which can cause
4421:     # division-by-zero in subsequent Bayes factor calculation.
4422:     prob_a_better = float(np.mean(samples_a > samples_b))
4423:     prob_a_better = float(np.clip(prob_a_better, 1e-6, ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L750_59", 1.0) - 1e-6))
4424:
4425:     # Compute Bayes factor (simplified)
4426:     if prob_a_better > ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L753_27", 0.5):
4427:         bayes_factor = prob_a_better / (1 - prob_a_better)
4428:     else:
4429:         bayes_factor = (1 - prob_a_better) / prob_a_better
4430:
4431:     return {
4432:         "probability_a_better": float(prob_a_better),
4433:         "bayes_factor": float(bayes_factor),
4434:         "difference_mean": float(np.mean(samples_a - samples_b)),
4435:         "difference_ci_95": (
4436:             float(
4437:                 np.percentile(
4438:                     samples_a - samples_b,
4439:                     2.5,
4440:                 )
4441:             ),
4442:             float(
4443:                 np.percentile(
4444:                     samples_a - samples_b,
4445:                     97.5,
4446:                 )
4447:             ),
4448:         ),
4449:     }
4450:
4451: # =====
4452: # CROSS-ENCODER RERANKING - State-of-the-Art Retrieval
4453: # =====
4454:
4455: class PolicyCrossEncoderReranker:
4456:     """
4457:     Cross-encoder reranking optimized for Spanish policy documents.
4458:
4459:     Uses transformer-based cross-attention for precise relevance scoring.
4460:     Superior to bi-encoder + cosine similarity for final ranking.
4461:     """
4462:
4463:     def __init__(
4464:         self,
4465:         model_name: str = DEFAULT_CROSS_ENCODER_MODEL,
4466:         max_length: int = 512,
4467:         retry_handler=None,
4468:     ) -> None:

```

```
4469:     """
4470:     Initialize cross-encoder reranker.
4471:
4472:     Args:
4473:         model_name: HuggingFace model name (multilingual preferred)
4474:         max_length: Maximum sequence length for cross-encoder
4475:         retry_handler: Optional RetryHandler for model loading
4476:
4477:     Raises:
4478:         RuntimeError: If online model download is required but HF_ONLINE=0
4479:     """
4480:     self._logger = logging.getLogger(self.__class__.__name__)
4481:     self.retry_handler = retry_handler
4482:
4483:     # Check dependency lockdown before attempting model load
4484:     from farfan_pipeline.core.dependency_lockdown import _is_model_cached, get_dependency_lockdown
4485:     lockdown = get_dependency_lockdown()
4486:
4487:     # Check if we're trying to download a remote model when offline
4488:     if not _is_model_cached(model_name):
4489:         lockdown.check_online_model_access(
4490:             model_name=model_name,
4491:             operation="load CrossEncoder model"
4492:         )
4493:
4494:     # Load model with retry logic if available
4495:     if retry_handler:
4496:         try:
4497:             from retry_handler import DependencyType
4498:
4499:             @retry_handler.with_retry(
4500:                 DependencyType.EMBEDDING_SERVICE,
4501:                 operation_name="load_cross_encoder",
4502:                 exceptions=(OSError, IOError, ConnectionError, RuntimeError)
4503:             )
4504:             def load_model():
4505:                 return CrossEncoder(model_name, max_length=max_length)
4506:
4507:                 self.model = load_model()
4508:                 self._logger.info(f"Cross-encoder loaded with retry protection: {model_name}")
4509:             except Exception as e:
4510:                 self._logger.error(f"Failed to load cross-encoder: {e}")
4511:                 raise
4512:         else:
4513:             self.model = CrossEncoder(model_name, max_length=max_length)
4514:             self._logger.info(f"Cross-encoder loaded: {model_name}")
4515:
4516:     def rerank(
4517:         self,
4518:         query: str,
4519:         candidates: list[SemanticChunk],
4520:         top_k: int = 10,
4521:         min_score: float = ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L848_27",
4522: 0.0),
4522:     ) -> list[tuple[SemanticChunk, float]]:
4523:         """
```

```

4524:     Rerank candidates using cross-encoder attention.
4525:
4526:     Returns top-k chunks with relevance scores.
4527:     """
4528:     if not candidates:
4529:         return []
4530:
4531:     # Prepare query-document pairs
4532:     pairs = [(query, chunk["content"]) for chunk in candidates]
4533:
4534:     # Score with cross-encoder
4535:     scores = self.model.predict(pairs, show_progress_bar=False)
4536:
4537:     # Combine chunks with scores and sort
4538:     ranked = sorted(zip(candidates, scores, strict=False), key=lambda x: x[1], reverse=True)
4539:
4540:     # Filter by minimum score and limit to top_k
4541:     filtered = [
4542:         (chunk, float(score)) for chunk, score in ranked if score >= min_score
4543:     ][:top_k]
4544:
4545:     self._logger.info(
4546:         "Reranked %d candidates, returned %d with min_score=%.2f",
4547:         len(candidates),
4548:         len(filtered),
4549:         min_score,
4550:     )
4551:
4552:     return filtered
4553:
4554: # =====
4555: # MAIN EMBEDDING SYSTEM - Orchestrator
4556: # =====
4557:
4558: @dataclass
4559: class PolicyEmbeddingConfig:
4560:     """Configuration for policy embedding system."""
4561:
4562:     # Model selection
4563:     embedding_model: str = MODEL_PARAPHRASE_MULTILINGUAL
4564:     cross_encoder_model: str = DEFAULT_CROSS_ENCODER_MODEL
4565:
4566:     # Chunking parameters
4567:     chunk_size: int = 512
4568:     chunk_overlap: int = 128
4569:
4570:     # Retrieval parameters
4571:     top_k_candidates: int = 50 # Bi-encoder retrieval
4572:     top_k_rerank: int = 10 # Cross-encoder rerank
4573:     mmr_lambda: float = ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L900_24", 0.
7) # Diversity vs relevance trade-off
4574:
4575:     # Bayesian analysis
4576:     prior_strength: float = ParameterLoaderV2.get("farfan_core.processing.embedding_policy.BayesianNumericalAnalyzer._null_evaluation", "auto_param_L903_28"
, 1.0) # Weakly informative prior
4577:

```

```
4578:     # Performance
4579:     batch_size: int = 32
4580:     normalize_embeddings: bool = True
4581:
4582: class PolicyAnalysisEmbedder:
4583:     """
4584:         Production-ready embedding system for Colombian PDM analysis.
4585:
4586:         Implements complete pipeline:
4587:             1. Advanced semantic chunking with P-D-Q awareness
4588:             2. Multilingual embedding (Spanish-optimized)
4589:             3. Bi-encoder retrieval + cross-encoder reranking
4590:             4. Bayesian numerical analysis with uncertainty quantification
4591:             5. MMR-based diversification
4592:
4593:         Thread-safe, production-grade, fully typed.
4594:     """
4595:
4596:     def __init__(self, config: PolicyEmbeddingConfig, retry_handler=None) -> None:
4597:         self.config = config
4598:         self._logger = logging.getLogger(self.__class__.__name__)
4599:         self.retry_handler = retry_handler
4600:
4601:         # Check dependency lockdown before attempting model loads
4602:         from farfan_pipeline.core.dependency_lockdown import _is_model_cached, get_dependency_lockdown
4603:         lockdown = get_dependency_lockdown()
4604:
4605:         # Check if we're trying to download remote models when offline
4606:         if not _is_model_cached(config.embedding_model):
4607:             lockdown.check_online_model_access(
4608:                 model_name=config.embedding_model,
4609:                 operation="load SentenceTransformer embedding model"
4610:             )
4611:
4612:         # Initialize embedding model with retry logic
4613:         if retry_handler:
4614:             try:
4615:                 from retry_handler import DependencyType
4616:
4617:                 @retry_handler.with_retry(
4618:                     DependencyType.EMBEDDING_SERVICE,
4619:                     operation_name="load_sentence_transformer",
4620:                     exceptions=(OSError, IOError, ConnectionError, RuntimeError)
4621:                 )
4622:                 def load_embedding_model():
4623:                     return SentenceTransformer(config.embedding_model)
4624:
4625:                     self._logger.info("Initializing embedding model with retry: %s", config.embedding_model)
4626:                     self.embedding_model = load_embedding_model()
4627:                 except Exception as e:
4628:                     self._logger.error(f"Failed to load embedding model: {e}")
4629:                     raise
4630:
4631:             else:
4632:                 self._logger.info("Initializing embedding model: %s", config.embedding_model)
4633:                 self.embedding_model = SentenceTransformer(config.embedding_model)
```

```
4634:     # Initialize cross-encoder with retry logic
4635:     self._logger.info("Initializing cross-encoder: %s", config.cross_encoder_model)
4636:     self.cross_encoder = PolicyCrossEncoderReranker(
4637:         config.cross_encoder_model,
4638:         retry_handler=retry_handler
4639:     )
4640:
4641:     self.chunker = AdvancedSemanticChunker(
4642:         ChunkingConfig(
4643:             chunk_size=config.chunk_size,
4644:             chunk_overlap=config.chunk_overlap,
4645:         )
4646:     )
4647:
4648:     self.bayesian_analyzer = BayesianNumericalAnalyzer(
4649:         prior_strength=config.prior_strength
4650:     )
4651:
4652:     # Cache
4653:     self._embedding_cache: dict[str, NDArray[np.float32]] = {}
4654:     self._chunk_cache: dict[str, list[SemanticChunk]] = {}
4655:
4656:     def process_document(
4657:         self,
4658:         document_text: str,
4659:         document_metadata: dict[str, Any],
4660:     ) -> list[SemanticChunk]:
4661:         """
4662:             Process complete PDM document into semantic chunks with embeddings.
4663:
4664:             Args:
4665:                 document_text: Full document text
4666:                 document_metadata: Metadata including doc_id, municipality, year
4667:
4668:             Returns:
4669:                 List of semantic chunks with embeddings and P-D-Q context
4670:         """
4671:         doc_id = document_metadata.get("doc_id", "unknown")
4672:         self._logger.info("Processing document: %s", doc_id)
4673:
4674:         # Check cache
4675:         if doc_id in self._chunk_cache:
4676:             self._logger.info(
4677:                 "Retrieved %d chunks from cache", len(self._chunk_cache[doc_id])
4678:             )
4679:             return self._chunk_cache[doc_id]
4680:
4681:         # Chunk document with semantic awareness
4682:         chunks = self.chunker.chunk_document(document_text, document_metadata)
4683:
4684:         # Generate embeddings in batches
4685:         chunk_texts = [chunk["content"] for chunk in chunks]
4686:         embeddings = self._embed_texts(chunk_texts)
4687:
4688:         # Attach embeddings to chunks
4689:         for chunk, embedding in zip(chunks, embeddings, strict=False):
```

```

4690:         chunk["embedding"] = embedding
4691:
4692:     # Cache results
4693:     self._chunk_cache[doc_id] = chunks
4694:
4695:     self._logger.info(
4696:         "Processed document %s: %d chunks, avg tokens: %.1f",
4697:         doc_id,
4698:         len(chunks),
4699:         np.mean([c["token_count"] for c in chunks]),
4700:     )
4701:
4702:     return chunks
4703:
4704: def semantic_search(
4705:     self,
4706:     query: str,
4707:     document_chunks: list[SemanticChunk],
4708:     pdq_filter: PDQIdentifier | None = None,
4709:     use_reranking: bool = True,
4710: ) -> list[tuple[SemanticChunk, float]]:
4711:     """
4712:         Advanced semantic search with P-D-Q filtering and reranking.
4713:
4714:         Pipeline:
4715:             1. Bi-encoder retrieval (fast, approximate)
4716:             2. P-D-Q filtering (if specified)
4717:             3. Cross-encoder reranking (precise)
4718:             4. MMR diversification
4719:
4720:         Args:
4721:             query: Search query
4722:             document_chunks: Pool of chunks to search
4723:             pdq_filter: Optional P-D-Q context filter
4724:             use_reranking: Enable cross-encoder reranking
4725:
4726:         Returns:
4727:             Ranked list of (chunk, score) tuples
4728:             """
4729:     if not document_chunks:
4730:         return []
4731:
4732:     # Bi-encoder retrieval: fast approximate search
4733:     chunk_embeddings = np.vstack([c["embedding"] for c in document_chunks])
4734:     query_embedding = self._embed_texts([query])[0]
4735:     similarities = cosine_similarity(
4736:         query_embedding.reshape(1, -1), chunk_embeddings
4737:     ).ravel()
4738:
4739:     # Get top-k candidates
4740:     top_indices = np.argsort(-similarities)[: self.config.top_k_candidates]
4741:     candidates = [document_chunks[i] for i in top_indices]
4742:
4743:     # Apply P-D-Q filter if specified
4744:     if pdq_filter:
4745:         candidates = self._filter_by_pdq(candidates, pdq_filter)

```

```
4746:             self._logger.info(
4747:                 "Filtered to %d chunks matching P-D-Q context", len(candidates)
4748:             )
4749:
4750:         if not candidates:
4751:             return []
4752:
4753:         # Cross-encoder reranking for precision
4754:         if use_reranking:
4755:             reranked = self.cross_encoder.rerank(
4756:                 query, candidates, top_k=self.config.top_k_rerank
4757:             )
4758:         else:
4759:             # Use bi-encoder scores
4760:             candidate_indices = [document_chunks.index(c) for c in candidates]
4761:             reranked = [
4762:                 (candidates[i], float(similarities[candidate_indices[i]]))
4763:                 for i in range(len(candidates))
4764:             ]
4765:             reranked.sort(key=lambda x: x[1], reverse=True)
4766:             reranked = reranked[: self.config.top_k_rerank]
4767:
4768:         # MMR diversification
4769:         if len(reranked) > 1:
4770:             reranked = self._apply_mmr(reranked)
4771:
4772:         return reranked
4773:
4774:     def evaluate_policy_numerical_consistency(
4775:         self,
4776:         chunks: list[SemanticChunk],
4777:         pdq_context: PDQIdentifier,
4778:     ) -> BayesianEvaluation:
4779:         """
4780:             Bayesian evaluation of numerical consistency for policy metric.
4781:
4782:             Extracts numerical values from chunks matching P-D-Q context,
4783:             performs rigorous statistical analysis with uncertainty quantification.
4784:
4785:             Args:
4786:                 chunks: Document chunks to analyze
4787:                 pdq_context: P-D-Q context to filter relevant chunks
4788:
4789:             Returns:
4790:                 Bayesian evaluation with credible intervals and evidence strength
4791:         """
4792:         # Filter chunks by P-D-Q context
4793:         relevant_chunks = self._filter_by_pdq(chunks, pdq_context)
4794:
4795:         if not relevant_chunks:
4796:             self._logger.warning(
4797:                 "No chunks found for P-D-Q context: %s",
4798:                 pdq_context["question_unique_id"],
4799:             )
4800:             return self.bayesian_analyzer._null_evaluation()
4801:
```

```
4802:     # Extract numerical values from chunks
4803:     numerical_values = self._extract_numerical_values(relevant_chunks)
4804:
4805:     if not numerical_values:
4806:         self._logger.warning(
4807:             "No numerical values extracted from %d chunks", len(relevant_chunks)
4808:         )
4809:     return self.bayesian_analyzer._null_evaluation()
4810:
4811:     # Perform Bayesian evaluation
4812:     evaluation = self.bayesian_analyzer.evaluate_policy_metric(numerical_values)
4813:
4814:     self._logger.info(
4815:         "Evaluated %d numerical values for %s: point_estimate=%.3f, CI=[%.3f, %.3f], evidence=%s",
4816:         len(numerical_values),
4817:         pdq_context["rubric_key"],
4818:         evaluation["point_estimate"],
4819:         evaluation["credible_interval_95"][0],
4820:         evaluation["credible_interval_95"][1],
4821:         evaluation["evidence_strength"],
4822:     )
4823:
4824:     return evaluation
4825:
4826: def compare_policy_interventions(
4827:     self,
4828:     intervention_a_chunks: list[SemanticChunk],
4829:     intervention_b_chunks: list[SemanticChunk],
4830:     pdq_context: PDQIdentifier,
4831: ) -> dict[str, Any]:
4832:     """
4833:     Bayesian comparison of two policy interventions.
4834:
4835:     Returns probability and evidence for superiority.
4836:     """
4837:     values_a = self._extract_numerical_values(
4838:         self._filter_by_pdq(intervention_a_chunks, pdq_context)
4839:     )
4840:     values_b = self._extract_numerical_values(
4841:         self._filter_by_pdq(intervention_b_chunks, pdq_context)
4842:     )
4843:
4844:     return self.bayesian_analyzer.compare_policies(values_a, values_b)
4845:
4846: def generate_pdq_report(
4847:     self,
4848:     document_chunks: list[SemanticChunk],
4849:     target_pdq: PDQIdentifier,
4850: ) -> dict[str, Any]:
4851:     """
4852:     Generate comprehensive analytical report for P-D-Q question.
4853:
4854:     Combines semantic search, numerical analysis, and evidence synthesis.
4855:     """
4856:     # Semantic search for relevant content
4857:     query = self._generate_query_from_pdq(target_pdq)
```

```
4858:         relevant_chunks = self.semantic_search(
4859:             query, document_chunks, pdq_filter=target_pdq
4860:         )
4861:
4862:         # Numerical consistency analysis
4863:         numerical_eval = self.evaluate_policy_numerical_consistency(
4864:             document_chunks, target_pdq
4865:         )
4866:
4867:         # Extract key evidence passages
4868:         evidence_passages = [
4869:             {
4870:                 "content": chunk["content"][:300],
4871:                 "relevance_score": float(score),
4872:                 "metadata": chunk["metadata"],
4873:             }
4874:             for chunk, score in relevant_chunks[:3]
4875:         ]
4876:
4877:         # Synthesize report
4878:         report = {
4879:             "question_unique_id": target_pdq["question_unique_id"],
4880:             "rubric_key": target_pdq["rubric_key"],
4881:             "evidence_count": len(relevant_chunks),
4882:             "numerical_evaluation": {
4883:                 "point_estimate": numerical_eval["point_estimate"],
4884:                 "credible_interval_95": numerical_eval["credible_interval_95"],
4885:                 "evidence_strength": numerical_eval["evidence_strength"],
4886:                 "numerical_coherence": numerical_eval["numerical_coherence"],
4887:             },
4888:             "evidence_passages": evidence_passages,
4889:             "confidence": self._compute_overall_confidence(
4890:                 relevant_chunks, numerical_eval
4891:             ),
4892:         }
4893:
4894:         return report
4895:
4896:     # =====
4897:     # PRIVATE METHODS
4898:     # =====
4899:
4900:     @calibrated_method("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._embed_texts")
4901:     def _embed_texts(self, texts: list[str]) -> NDArray[np.float32]:
4902:         """Generate embeddings with caching and retry logic."""
4903:         uncached_texts = []
4904:         uncached_indices = []
4905:
4906:         embeddings_list = []
4907:
4908:         for i, text in enumerate(texts):
4909:             text_hash = hashlib.sha256(text.encode()).hexdigest()[:16]
4910:
4911:             if text_hash in self._embedding_cache:
4912:                 embeddings_list.append(self._embedding_cache[text_hash])
4913:             else:
```

```
4914:         uncached_texts.append(text)
4915:         uncached_indices.append((i, text_hash))
4916:         embeddings_list.append(None) # Placeholder
4917:
4918:     # Generate embeddings for uncached texts with retry logic
4919:     if uncached_texts:
4920:         if self.retry_handler:
4921:             try:
4922:                 from retry_handler import DependencyType
4923:
4924:                 @self.retry_handler.with_retry(
4925:                     DependencyType.EMBEDDING_SERVICE,
4926:                     operation_name="encode_texts",
4927:                     exceptions=(ConnectionError, TimeoutError, RuntimeError, OSError)
4928:                 )
4929:             def encode_with_retry():
4930:                 return self.embedding_model.encode(
4931:                     uncached_texts,
4932:                     batch_size=self.config.batch_size,
4933:                     normalize_embeddings=self.config.normalize_embeddings,
4934:                     show_progress_bar=False,
4935:                     convert_to_numpy=True,
4936:                 )
4937:
4938:                 new_embeddings = encode_with_retry()
4939:             except Exception as e:
4940:                 self._logger.error(f"Failed to encode texts with retry: {e}")
4941:                 raise
4942:             else:
4943:                 new_embeddings = self.embedding_model.encode(
4944:                     uncached_texts,
4945:                     batch_size=self.config.batch_size,
4946:                     normalize_embeddings=self.config.normalize_embeddings,
4947:                     show_progress_bar=False,
4948:                     convert_to_numpy=True,
4949:                 )
4950:
4951:             # Cache and insert
4952:             for (orig_idx, text_hash), emb in zip(uncached_indices, new_embeddings, strict=False):
4953:                 self._embedding_cache[text_hash] = emb
4954:                 embeddings_list[orig_idx] = emb
4955:
4956:             return np.vstack(embeddings_list).astype(np.float32)
4957:
4958:     def _filter_by_pdq(
4959:         self, chunks: list[SemanticChunk], pdq_filter: PDQIdentifier
4960:     ) -> list[SemanticChunk]:
4961:         """Filter chunks by P-D-Q context."""
4962:
4963:     def __repr__(value: Any) -> str:
4964:         if value is None or isinstance(value, (int, float, bool)):
4965:             return repr(value)
4966:         if isinstance(value, str):
4967:             # Strip excessive whitespace for logging clarity
4968:             preview = value if len(value) <= 24 else f"{value[:21]}..."
4969:             return repr(preview)
```

```
4970:         return type(value).__name__
4971:
4972:     def _log_mismatch(key: str, needed: Any, index: int | None = None) -> None:
4973:         message = (
4974:             "ERR_CONTRACT_MISMATCH[fn=_filter_by_pdq, "
4975:             f"key='{key}', needed={_repr_contract(needed)}, got={_repr_contract(got)}]"
4976:         )
4977:         if index is not None:
4978:             message += f", index={index}"
4979:         message += "]"
4980:         self._logger.error(message)
4981:
4982:         self._logger.debug(
4983:             "edge %s \u206\222 _filter_by_pdq | params=%s",
4984:             self.__class__.__name__,
4985:             {
4986:                 "chunks_type": type(chunks).__name__,
4987:                 "chunks_len": len(chunks) if isinstance(chunks, list) else "n/a",
4988:                 "pdq_filter_type": type(pdq_filter).__name__,
4989:                 "pdq_filter_keys": sorted(pdq_filter.keys())
4990:                 if isinstance(pdq_filter, dict)
4991:                 else None,
4992:             },
4993:         )
4994:
4995:         if not isinstance(chunks, list):
4996:             _log_mismatch("chunks", "list", chunks)
4997:         return []
4998:
4999:         if not isinstance(pdq_filter, dict):
5000:             _log_mismatch("pdq_filter", "dict", pdq_filter)
5001:         return []
5002:
5003:         expected_policy = pdq_filter.get("policy")
5004:         expected_dimension = pdq_filter.get("dimension")
5005:
5006:         if expected_policy is None or expected_dimension is None:
5007:             _log_mismatch(
5008:                 "pdq_filter",
5009:                 "keys='('policy','dimension')",
5010:                 {"policy": expected_policy, "dimension": expected_dimension},
5011:             )
5012:         return []
5013:
5014:         filtered_chunks: list[SemanticChunk] = []
5015:
5016:         for index, chunk in enumerate(chunks):
5017:             if not isinstance(chunk, dict):
5018:                 _log_mismatch("chunk", "dict", chunk, index)
5019:                 continue
5020:
5021:             pdq_context = chunk.get("pdq_context")
5022:
5023:             if not pdq_context:
5024:                 _log_mismatch("pdq_context", True, pdq_context, index)
5025:                 continue
```

```
5026:
5027:         if not isinstance(pdq_context, dict):
5028:             _log_mismatch("pdq_context", "dict", pdq_context, index)
5029:             continue
5030:
5031:         policy = pdq_context.get("policy")
5032:         dimension = pdq_context.get("dimension")
5033:
5034:         if policy is None or dimension is None:
5035:             _log_mismatch(
5036:                 "pdq_context",
5037:                 "keys=('policy','dimension')",
5038:                 {"policy": policy, "dimension": dimension},
5039:                 index,
5040:             )
5041:             continue
5042:
5043:         if policy == expected_policy and dimension == expected_dimension:
5044:             filtered_chunks.append(chunk)
5045:
5046:     return filtered_chunks
5047:
5048: def _apply_mmr(
5049:     self,
5050:     ranked_results: list[tuple[SemanticChunk, float]],
5051: ) -> list[tuple[SemanticChunk, float]]:
5052:     """
5053:     Apply Maximal Marginal Relevance for diversification.
5054:
5055:     Balances relevance with diversity to avoid redundant results.
5056:     """
5057:     if len(ranked_results) <= 1:
5058:         return ranked_results
5059:
5060:     chunks, scores = zip(*ranked_results, strict=False)
5061:     chunk_embeddings = np.vstack([c["embedding"] for c in chunks])
5062:
5063:     selected_indices = []
5064:     remaining_indices = list(range(len(chunks)))
5065:
5066:     # Select first (most relevant)
5067:     selected_indices.append(0)
5068:     remaining_indices.remove(0)
5069:
5070:     # Iteratively select diverse documents
5071:     while remaining_indices and len(selected_indices) < len(chunks):
5072:         best_mmr_score = float("-inf")
5073:         best_idx = None
5074:
5075:         for idx in remaining_indices:
5076:             # Relevance score
5077:             relevance = scores[idx]
5078:
5079:             # Diversity: max similarity to selected
5080:             similarities_to_selected = cosine_similarity(
5081:                 chunk_embeddings[idx : idx + 1],
```



```

5138:                 num_str = raw_num.replace(", ", ".")
5139:             else:
5140:                 # Only dot or plain number
5141:                 num_str = raw_num
5142:
5143:                 value = float(num_str)
5144:
5145:                 # Normalize to 0-1 scale if it's a percentage
5146:                 if "%" in match.group(0) and value <= 100:
5147:                     value = value / 100.0
5148:
5149:                 # Filter outliers
5150:                 if 0 <= value <= 1e9: # Reasonable range
5151:                     numerical_values.append(value)
5152:
5153:             except (ValueError, IndexError):
5154:                 continue
5155:
5156:         return numerical_values
5157:
5158: @calibrated_method("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq")
5159: def _generate_query_from_pdq(self, pdq: PDQIdentifier) -> str:
5160:     """Generate search query from P-D-Q identifier."""
5161:     policy_name = PolicyDomain[pdq["policy"]].value
5162:     dimension_name = AnalyticalDimension[pdq["dimension"]].value
5163:
5164:     query = f"{policy_name} - {dimension_name}"
5165:     return query
5166:
5167: def _compute_overall_confidence(
5168:     self,
5169:     relevant_chunks: list[tuple[SemanticChunk, float]],
5170:     numerical_eval: BayesianEvaluation,
5171: ) -> float:
5172:     """
5173:         Compute overall confidence score combining semantic and numerical evidence.
5174:
5175:     Considerations:
5176:     - Number of relevant chunks
5177:     - Semantic relevance scores
5178:     - Numerical evidence strength
5179:     - Statistical coherence
5180:     """
5181:     if not relevant_chunks:
5182:         return ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L1509_19", 0)
5183:
5184:     # Semantic confidence: average of top scores
5185:     semantic_scores = [score for _, score in relevant_chunks[:5]]
5186:     semantic_confidence = (
5187:         float(np.mean(semantic_scores)) if semantic_scores else ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._g
5188:         generate_query_from_pdq", "auto_param_L1514_68", 0.0)
5189:     )
5190:
5191:     # Numerical confidence: based on evidence strength and coherence

```

```

5192:         "weak": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L1519_20",
0.25),
5193:         "moderate": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L1520_2
4", 0.5),
5194:         "strong": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L1521_22"
, 0.75),
5195:         "very_strong": ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L152
2_27", 1.0),
5196:     }
5197:     numerical_confidence = (
5198:         evidence_strength_map[numerical_eval["evidence_strength"]]
5199:         * numerical_eval["numerical_coherence"]
5200:     )
5201:
5202:     # Combined confidence: weighted average
5203:     overall_confidence = ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L1
530_29", 0.6) * semantic_confidence + ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_
L1530_57", 0.4) * numerical_confidence
5204:
5205:     return float(np.clip(overall_confidence, ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_
pdq", "auto_param_L1532_49", 0.0), ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._generate_query_from_pdq", "auto_param_L15
32_54", 1.0)))
5206:
5207:     @lru_cache(maxsize=1024)
5208:     @calibrated_method("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder._cached_similarity")
5209:     def _cached_similarity(self, text_hash1: str, text_hash2: str) -> float:
5210:         """Cached similarity computation for performance.
5211:         Assumes embeddings are cached in self._embedding_cache using text_hash as key.
5212:         """
5213:         emb1 = self._embedding_cache[text_hash1]
5214:         emb2 = self._embedding_cache[text_hash2]
5215:         return float(cosine_similarity(emb1.reshape(1, -1), emb2.reshape(1, -1))[0, 0])
5216:
5217:     @calibrated_method("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder.get_diagnostics")
5218:     def get_diagnostics(self) -> dict[str, Any]:
5219:         """Get system diagnostics and performance metrics."""
5220:         return {
5221:             "model": self.config.embedding_model,
5222:             "embedding_cache_size": len(self._embedding_cache),
5223:             "chunk_cache_size": len(self._chunk_cache),
5224:             "total_chunks_processed": sum(
5225:                 len(chunks) for chunks in self._chunk_cache.values()
5226:             ),
5227:             "config": {
5228:                 "chunk_size": self.config.chunk_size,
5229:                 "chunk_overlap": self.config.chunk_overlap,
5230:                 "top_k_candidates": self.config.top_k_candidates,
5231:                 "top_k_rerank": self.config.top_k_rerank,
5232:                 "mmr_lambda": self.config.mmr_lambda,
5233:             },
5234:         }
5235:
5236:     # =====
5237:     # PRODUCTION FACTORY AND UTILITIES
5238:     # =====
5239:

```

```
5240: def create_policy_embedder(
5241:     model_tier: Literal["fast", "balanced", "accurate"] = "balanced",
5242: ) -> PolicyAnalysisEmbedder:
5243:     """
5244:         Factory function for creating production-ready policy embedder.
5245:
5246:     Args:
5247:         model_tier: Performance/accuracy trade-off
5248:             - "fast": Lightweight, low latency
5249:             - "balanced": Good performance/accuracy balance (default)
5250:             - "accurate": Maximum accuracy, higher latency
5251:
5252:     Returns:
5253:         Configured PolicyAnalysisEmbedder instance
5254:     """
5255:     model_configs = {
5256:         "fast": PolicyEmbeddingConfig(
5257:             embedding_model="sentence-transformers/paraphrase-multilingual-MiniLM-L12-v2",
5258:             cross_encoder_model=DEFAULT_CROSS_ENCODER_MODEL,
5259:             chunk_size=256,
5260:             chunk_overlap=64,
5261:             top_k_candidates=30,
5262:             top_k_rerank=5,
5263:             batch_size=64,
5264:         ),
5265:         "balanced": PolicyEmbeddingConfig(
5266:             embedding_model=MODEL_PARAPHRASE_MULTILINGUAL,
5267:             cross_encoder_model=DEFAULT_CROSS_ENCODER_MODEL,
5268:             chunk_size=512,
5269:             chunk_overlap=128,
5270:             top_k_candidates=50,
5271:             top_k_rerank=10,
5272:             batch_size=32,
5273:         ),
5274:         "accurate": PolicyEmbeddingConfig(
5275:             embedding_model=MODEL_PARAPHRASE_MULTILINGUAL,
5276:             cross_encoder_model="cross-encoder/mmarco-mMiniLMv2-L12-H384-v1",
5277:             chunk_size=768,
5278:             chunk_overlap=192,
5279:             top_k_candidates=100,
5280:             top_k_rerank=20,
5281:             batch_size=16,
5282:         ),
5283:     }
5284:
5285:     config = model_configs[model_tier]
5286:
5287:     logger = logging.getLogger("PolicyEmbedderFactory")
5288:     logger.info("Creating policy embedder with tier: %s", model_tier)
5289:
5290:     return PolicyAnalysisEmbedder(config)
5291:
5292: # =====
5293: # PRODUCER CLASS - Registry Exposure
5294: # =====
5295:
```

```
5296: class EmbeddingPolicyProducer:
5297:     """
5298:         Producer wrapper for embedding policy analysis with registry exposure
5299:
5300:         Provides public API methods for orchestrator integration without exposing
5301:             internal implementation details or summarization logic.
5302:
5303:         Version: ParameterLoaderV2.get("farfan_core.processing.embedding_policy.PolicyAnalysisEmbedder.get_diagnostics", "auto_param_L1630_13", 1.0).0
5304:         Producer Type: Embedding / Semantic Search
5305:     """
5306:
5307:     def __init__(
5308:         self,
5309:             config: PolicyEmbeddingConfig | None = None,
5310:             model_tier: Literal["fast", "balanced", "accurate"] = "balanced",
5311:             retry_handler=None
5312:     ) -> None:
5313:         """Initialize producer with optional configuration"""
5314:         if config is None:
5315:             self.embedder = create_policy_embedder(model_tier)
5316:         else:
5317:             self.embedder = PolicyAnalysisEmbedder(config, retry_handler=retry_handler)
5318:
5319:         self._logger = logging.getLogger(self.__class__.__name__)
5320:         self._logger.info("EmbeddingPolicyProducer initialized")
5321:
5322: # =====
5323: # DOCUMENT PROCESSING API
5324: # =====
5325:
5326:     def process_document(
5327:         self,
5328:             document_text: str,
5329:             document_metadata: dict[str, Any]
5330:     ) -> list[SemanticChunk]:
5331:         """Process document into semantic chunks with embeddings"""
5332:         return self.embedder.process_document(document_text, document_metadata)
5333:
5334:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_chunk_count")
5335:     def get_chunk_count(self, chunks: list[SemanticChunk]) -> int:
5336:         """Get number of chunks"""
5337:         return len(chunks)
5338:
5339:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_chunk_text")
5340:     def get_chunk_text(self, chunk: SemanticChunk) -> str:
5341:         """Extract text from chunk"""
5342:         return chunk["content"]
5343:
5344:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_chunk_embedding")
5345:     def get_chunk_embedding(self, chunk: SemanticChunk) -> NDArray[np.float32]:
5346:         """Extract embedding from chunk"""
5347:         return chunk["embedding"]
5348:
5349:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_chunk_metadata")
5350:     def get_chunk_metadata(self, chunk: SemanticChunk) -> dict[str, Any]:
5351:         """Extract metadata from chunk"""
```

```
5352:         return chunk["metadata"]
5353:
5354:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_chunk_pdq_context")
5355:     def get_chunk_pdq_context(self, chunk: SemanticChunk) -> PDQIdentifier | None:
5356:         """Extract P-D-Q context from chunk"""
5357:         return chunk["pdq_context"]
5358:
5359:     # =====
5360:     # SEMANTIC SEARCH API
5361:     # =====
5362:
5363:     def semantic_search(
5364:         self,
5365:         query: str,
5366:         document_chunks: list[SemanticChunk],
5367:         pdq_filter: PDQIdentifier | None = None,
5368:         use_reranking: bool = True
5369:     ) -> list[tuple[SemanticChunk, float]]:
5370:         """Advanced semantic search with reranking"""
5371:         return self.embedder.semantic_search(
5372:             query, document_chunks, pdq_filter, use_reranking
5373:         )
5374:
5375:     def get_search_result_chunk(
5376:         self, result: tuple[SemanticChunk, float]
5377:     ) -> SemanticChunk:
5378:         """Extract chunk from search result"""
5379:         return result[0]
5380:
5381:     def get_search_result_score(
5382:         self, result: tuple[SemanticChunk, float]
5383:     ) -> float:
5384:         """Extract relevance score from search result"""
5385:         return result[1]
5386:
5387:     # =====
5388:     # P-D-Q ANALYSIS API
5389:     # =====
5390:
5391:     def generate_pdq_report(
5392:         self,
5393:         document_chunks: list[SemanticChunk],
5394:         target_pdq: PDQIdentifier
5395:     ) -> dict[str, Any]:
5396:         """Generate comprehensive analytical report for P-D-Q question"""
5397:         return self.embedder.generate_pdq_report(document_chunks, target_pdq)
5398:
5399:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_pdq_evidence_count")
5400:     def get_pdq_evidence_count(self, report: dict[str, Any]) -> int:
5401:         """Extract evidence count from P-D-Q report"""
5402:         return report.get("evidence_count", 0)
5403:
5404:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_pdq_numerical_evaluation")
5405:     def get_pdq_numerical_evaluation(self, report: dict[str, Any]) -> dict[str, Any]:
5406:         """Extract numerical evaluation from P-D-Q report"""
5407:         return report.get("numerical_evaluation", {})
```

```
5408:
5409:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_pdq_evidence_passages")
5410:     def get_pdq_evidence_passages(self, report: dict[str, Any]) -> list[dict[str, Any]]:
5411:         """Extract evidence passages from P-D-Q report"""
5412:         return report.get("evidence_passages", [])
5413:
5414:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_pdq_confidence")
5415:     def get_pdq_confidence(self, report: dict[str, Any]) -> float:
5416:         """Extract confidence from P-D-Q report"""
5417:         return report.get("confidence", ParameterLoaderV2.get("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_pdq_confidence", "auto_param_L1744_40", 0.0))
5418:
5419:     # =====
5420:     # BAYESIAN NUMERICAL ANALYSIS API
5421:     # =====
5422:
5423:     def evaluate_numerical_consistency(
5424:         self,
5425:         chunks: list[SemanticChunk],
5426:         pdq_context: PDQIdentifier
5427:     ) -> BayesianEvaluation:
5428:         """Evaluate numerical consistency with Bayesian analysis"""
5429:         return self.embedder.evaluate_policy_numerical_consistency(
5430:             chunks, pdq_context
5431:         )
5432:
5433:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_point_estimate")
5434:     def get_point_estimate(self, evaluation: BayesianEvaluation) -> float:
5435:         """Extract point estimate from Bayesian evaluation"""
5436:         return evaluation["point_estimate"]
5437:
5438:     def get_credible_interval(
5439:         self, evaluation: BayesianEvaluation
5440:     ) -> tuple[float, float]:
5441:         """Extract 95% credible interval from Bayesian evaluation"""
5442:         return evaluation["credible_interval_95"]
5443:
5444:     def get_evidence_strength(
5445:         self, evaluation: BayesianEvaluation
5446:     ) -> Literal["weak", "moderate", "strong", "very_strong"]:
5447:         """Extract evidence strength classification"""
5448:         return evaluation["evidence_strength"]
5449:
5450:     @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_numerical_coherence")
5451:     def get_numerical_coherence(self, evaluation: BayesianEvaluation) -> float:
5452:         """Extract numerical coherence score"""
5453:         return evaluation["numerical_coherence"]
5454:
5455:     # =====
5456:     # POLICY COMPARISON API
5457:     # =====
5458:
5459:     def compare_policy_interventions(
5460:         self,
5461:         intervention_a_chunks: list[SemanticChunk],
5462:         intervention_b_chunks: list[SemanticChunk],
```

```
5463:     pdq_context: PDQIdentifier
5464: ) -> dict[str, Any]:
5465:     """Bayesian comparison of two policy interventions"""
5466:     return self.embedder.compare_policy_interventions(
5467:         intervention_a_chunks, intervention_b_chunks, pdq_context
5468:     )
5469:
5470: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_probability")
5471: def get_comparison_probability(self, comparison: dict[str, Any]) -> float:
5472:     """Extract probability that A is better than B"""
5473:     return comparison.get("probability_a_better", ParameterLoaderV2.get("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_probability", "auto_param_L1800_54", 0.5))
5474:
5475: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_bayes_factor")
5476: def get_comparison_bayes_factor(self, comparison: dict[str, Any]) -> float:
5477:     """Extract Bayes factor from comparison"""
5478:     return comparison.get("bayes_factor", ParameterLoaderV2.get("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_bayes_factor", "auto_param_L1805_46", 1.0))
5479:
5480: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_difference_mean")
5481: def get_comparison_difference_mean(self, comparison: dict[str, Any]) -> float:
5482:     """Extract mean difference from comparison"""
5483:     return comparison.get("difference_mean", ParameterLoaderV2.get("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_comparison_difference_mean", "auto_param_L1810_49", 0.0))
5484:
5485: # =====
5486: # UTILITY API
5487: # =====
5488:
5489: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_diagnostics")
5490: def get_diagnostics(self) -> dict[str, Any]:
5491:     """Get system diagnostics and performance metrics"""
5492:     return self.embedder.get_diagnostics()
5493:
5494: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_config")
5495: def get_config(self) -> PolicyEmbeddingConfig:
5496:     """Get current configuration"""
5497:     return self.embedder.config
5498:
5499: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.list_policy_domains")
5500: def list_policy_domains(self) -> list[PolicyDomain]:
5501:     """List all policy domains"""
5502:     return list(PolicyDomain)
5503:
5504: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.list_analytical_dimensions")
5505: def list_analytical_dimensions(self) -> list[AnalyticalDimension]:
5506:     """List all analytical dimensions"""
5507:     return list(AnalyticalDimension)
5508:
5509: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_policy_domain_description")
5510: def get_policy_domain_description(self, domain: PolicyDomain) -> str:
5511:     """Get description for policy domain"""
5512:     return domain.value
5513:
5514: @calibrated_method("farfan_core.processing.embedding_policy.EmbeddingPolicyProducer.get_analytical_dimension_description")
5515: def get_analytical_dimension_description(self, dimension: AnalyticalDimension) -> str:
```

```
5516:     """Get description for analytical dimension"""
5517:     return dimension.value
5518:
5519:     def create_pdq_identifier(
5520:         self,
5521:         policy: str,
5522:         dimension: str,
5523:         question: int
5524:     ) -> PDQIdentifier:
5525:         """Create P-D-Q identifier"""
5526:         return PDQIdentifier(
5527:             question_unique_id=f"{policy}-{dimension}-Q{question}",
5528:             policy=policy,
5529:             dimension=dimension,
5530:             question=question,
5531:             rubric_key=f"{dimension}-Q{question}"
5532:         )
5533:
5534: # =====
5535: # COMPREHENSIVE EXAMPLE - Production Usage
5536: # =====
5537:
5538: def example_pdm_analysis() -> None:
5539:     """
5540:     Complete example: analyzing Colombian Municipal Development Plan.
5541:     """
5542:     import logging
5543:
5544:     logging.basicConfig(level=logging.INFO)
5545:
5546:     # Sample PDM excerpt (simplified)
5547:     pdm_document = """
5548:     PLAN DE DESARROLLO MUNICIPAL 2024-2027
5549:     MUNICIPIO DE EJEMPLO, COLOMBIA
5550:
5551:     EJE ESTRATÃ\211GICO 1: DERECHOS DE LAS MUJERES E IGUALDAD DE GÃ\211NERO
5552:
5553:     DIAGNÃ\223STICO
5554:     El municipio presenta una brecha de gÃ@nero del 18.5% en participaciÃ³n laboral.
5555:     Se identificaron 2,340 mujeres en situaciÃ³n de vulnerabilidad econÃ³mica.
5556:     El presupuesto asignado asciende a $450 millones para el cuatrienio.
5557:
5558:     DISEÃ\221O DE INTERVENCIÃ\223N
5559:     Se implementarÃ¡n 3 programas de empoderamiento econÃ³mico:
5560:     - Programa de formaciÃ³n tÃ@cnica: 500 beneficiarias
5561:     - MicrocrÃ©ditos productivos: $280 millones
5562:     - Fortalecimiento empresarial: 150 emprendimientos
5563:
5564:     PRODUCTOS Y OUTPUTS
5565:     Meta cuatrienio: reducir brecha de gÃ@nero al 12% (reducciÃ³n del 35.1%)
5566:     Indicador: Tasa de participaciÃ³n laboral femenina
5567:     LÃ-nea base: 42.3% | Meta: 55.8%
5568:
5569:     RESULTADOS ESPERADOS
5570:     Incremento del 25% en ingresos promedio de beneficiarias
5571:     CreaciÃ³n de 320 nuevos empleos formales para mujeres
```

```
5572: Sostenibilidad: 78% de emprendimientos activos a 2 años
5573: """
5574:
5575:     metadata = {
5576:         "doc_id": "PDM_EJEMPLO_2024_2027",
5577:         "municipality": "Ejemplo",
5578:         "department": "Ejemplo",
5579:         "year": 2024,
5580:     }
5581:
5582:     # Create embedder
5583:     print("=" * 80)
5584:     print("POLICY ANALYSIS EMBEDDER - PRODUCTION EXAMPLE")
5585:     print("=" * 80)
5586:
5587:     embedder = create_policy_embedder(model_tier="balanced")
5588:
5589:     # Process document
5590:     print("\n1. PROCESSING DOCUMENT")
5591:     chunks = embedder.process_document(pdm_document, metadata)
5592:     print(f"    Generated {len(chunks)} semantic chunks")
5593:
5594:     # Define P-D-Q query
5595:     pdq_query = PDQIdentifier(
5596:         question_unique_id="P1-D1-Q3",
5597:         policy="P1",
5598:         dimension="D1",
5599:         question=3,
5600:         rubric_key="D1-Q3",
5601:     )
5602:
5603:     print(f"\n2. ANALYZING P-D-Q: {pdq_query['question_unique_id']}")
5604:     print(f"    Policy: {PolicyDomain.P1.value}")
5605:     print(f"    Dimension: {AnalyticalDimension.D1.value}")
5606:
5607:     # Generate comprehensive report
5608:     report = embedder.generate_pdq_report(chunks, pdq_query)
5609:
5610:     print("\n3. ANALYSIS RESULTS")
5611:     print(f"    Evidence chunks found: {report['evidence_count']}")
5612:     print(f"    Overall confidence: {report['confidence']:.3f}")
5613:     print("\n    Numerical Evaluation:")
5614:     print(
5615:         f"        - Point estimate: {report['numerical_evaluation']['point_estimate']:.3f}"
5616:     )
5617:     print(
5618:         f"        - 95% CI: [{report['numerical_evaluation']['credible_interval_95'][0]:.3f}, "
5619:             f"{report['numerical_evaluation']['credible_interval_95'][1]:.3f}]"
5620:     )
5621:     print(
5622:         f"        - Evidence strength: {report['numerical_evaluation']['evidence_strength']}"
5623:     )
5624:     print(
5625:         f"        - Numerical coherence: {report['numerical_evaluation']['numerical_coherence']:.3f}"
5626:     )
5627:
```

```
5628:     print("\n4. TOP EVIDENCE PASSAGES:")
5629:     for i, passage in enumerate(report["evidence_passages"], 1):
5630:         print(f"\n    [{i}] Relevance: {passage['relevance_score']:.3f}")
5631:         print(f"        {passage['content'][:200]}...")
5632:
5633:     # System diagnostics
5634:     print("\n5. SYSTEM DIAGNOSTICS")
5635:     diag = embedder.get_diagnostics()
5636:     print(f"    Model: {diag['model']}")
5637:     print(f"    Cache efficiency: {diag['embedding_cache_size']} embeddings cached")
5638:     print(f"    Total chunks processed: {diag['total_chunks_processed']}")
5639:
5640:     print("\n" + "=" * 80)
5641:     print("ANALYSIS COMPLETE")
5642:     print("=" * 80)
5643:
5644:
5645:
5646: =====
5647: FILE: src/farfan_pipeline/processing/models.py
5648: =====
5649:
5650: """
5651: CPP Ingestion Models (Deprecated - Use SPC)
5652:
5653: Data models for Canon Policy Package (CPP) ingestion pipeline.
5654: These models define the structure of policy documents after phase-one ingestion.
5655:
5656: NOTE: This is a compatibility layer. New code should use SPC (Smart Policy Chunks) terminology.
5657: """
5658:
5659: from __future__ import annotations
5660:
5661: from dataclasses import dataclass, field
5662: from enum import Enum
5663: from typing import Any
5664: from farfan_pipeline.core.parameters import ParameterLoaderV2
5665: from farfan_pipeline.core.calibration.decorators import calibrated_method
5666:
5667:
5668: class ChunkResolution(Enum):
5669:     """Granularity level for policy chunks."""
5670:     MICRO = "MICRO"  # Fine-grained chunks (sentences, clauses)
5671:     MESO = "MESO"    # Medium chunks (paragraphs, sections)
5672:     MACRO = "MACRO"  # Coarse chunks (chapters, themes)
5673:
5674:
5675: @dataclass
5676: class TextSpan:
5677:     """Represents a span of text in the original document."""
5678:     start: int
5679:     end: int
5680:
5681:
5682: @dataclass
5683: class Confidence:
```

```
5684:     """Confidence scores for various extraction processes."""
5685:     layout: float = 1.0
5686:     ocr: float = 1.0
5687:     typing: float = 1.0
5688:
5689:
5690: @dataclass
5691: class PolicyFacet:
5692:     """Policy-related metadata facets."""
5693:     programs: list[str] = field(default_factory=list)
5694:     projects: list[str] = field(default_factory=list)
5695:     axes: list[str] = field(default_factory=list)
5696:
5697:
5698: @dataclass
5699: class TimeFacet:
5700:     """Temporal metadata facets."""
5701:     years: list[int] = field(default_factory=list)
5702:     periods: list[str] = field(default_factory=list)
5703:
5704:
5705: @dataclass
5706: class GeoFacet:
5707:     """Geographic metadata facets."""
5708:     territories: list[str] = field(default_factory=list)
5709:     regions: list[str] = field(default_factory=list)
5710:
5711:
5712: @dataclass
5713: class ProvenanceMap:
5714:     """Provenance information for chunk extraction."""
5715:     source_page: int | None = None
5716:     source_section: str | None = None
5717:     extraction_method: str = "semantic_chunking"
5718:
5719:
5720: @dataclass
5721: class Budget:
5722:     """Budget information extracted from policy document."""
5723:     source: str
5724:     use: str
5725:     amount: float
5726:     year: int
5727:     currency: str = "COP"
5728:
5729:
5730: # Alias for compatibility
5731: BudgetInfo = Budget
5732:
5733:
5734: @dataclass
5735: class KPI:
5736:     """Key Performance Indicator extracted from policy."""
5737:     indicator_name: str
5738:     target_value: float | None = None
5739:     unit: str | None = None
```

```
5740:     year: int | None = None
5741:
5742:
5743: @dataclass
5744: class Entity:
5745:     """Named entity extracted from text."""
5746:     text: str
5747:     entity_type: str
5748:     confidence: float = 1.0
5749:
5750:
5751: @dataclass
5752: class Chunk:
5753:     """
5754:     A semantic chunk of policy text with metadata.
5755:
5756:     This is the fundamental unit of the CPP/SPC ingestion pipeline.
5757:     """
5758:     id: str
5759:     text: str
5760:     text_span: TextSpan
5761:     resolution: ChunkResolution
5762:     bytes_hash: str
5763:     policy_area_id: str | None = None # PA01-PA10 canonical code
5764:     dimension_id: str | None = None # DIM01-DIM06 canonical code
5765:
5766:     # Facets and metadata
5767:     policy_facets: PolicyFacet = field(default_factory=PolicyFacet)
5768:     time_facets: TimeFacet = field(default_factory=TimeFacet)
5769:     geo_facets: GeoFacet = field(default_factory=GeoFacet)
5770:     confidence: Confidence = field(default_factory=Confidence)
5771:
5772:     # Optional structured data
5773:     provenance: ProvenanceMap | None = None
5774:     budget: Budget | None = None
5775:     kpi: KPI | None = None
5776:     entities: list[Entity] = field(default_factory=list)
5777:
5778:
5779: @dataclass
5780: class ChunkGraph:
5781:     """
5782:     Graph structure containing all chunks and their relationships.
5783:     """
5784:     chunks: dict[str, Chunk] = field(default_factory=dict)
5785:     edges: list[tuple[str, str, str]] = field(default_factory=list) # (from_id, to_id, relation_type)
5786:
5787:     @calibrated_method("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_chunk")
5788:     def add_chunk(self, chunk: Chunk) -> None:
5789:         """Add a chunk to the graph."""
5790:         self.chunks[chunk.id] = chunk
5791:
5792:     @calibrated_method("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge")
5793:     def add_edge(self, from_id: str, to_id: str, relation_type: str) -> None:
5794:         """Add an edge to the graph."""
5795:         self.edges.append((from_id, to_id, relation_type))
```

```
5796:  
5797:  
5798: @dataclass  
5799: class PolicyManifest:  
5800:     """  
5801:         High-level manifest summarizing policy structure.  
5802:     """  
5803:     axes: list[str] = field(default_factory=list)  
5804:     programs: list[str] = field(default_factory=list)  
5805:     projects: list[str] = field(default_factory=list)  
5806:     years: list[int] = field(default_factory=list)  
5807:     territories: list[str] = field(default_factory=list)  
5808:     indicators: list[str] = field(default_factory=list)  
5809:     budget_rows: int = 0  
5810:  
5811:  
5812: @dataclass  
5813: class QualityMetrics:  
5814:     """  
5815:         Quality metrics for the ingestion process.  
5816:     """  
5817:     boundary_f1: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L167_25", 0.0)  
5818:     kpi_linkage_rate: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L168_30", 0.0)  
5819:     budget_consistency_score: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L169_38", 0.0)  
5820:     provenance_completeness: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L170_37", 0.0)  
5821:     structural_consistency: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L171_36", 0.0)  
5822:     temporal_robustness: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L172_33", 0.0)  
5823:     chunk_context_coverage: float = ParameterLoaderV2.get("farfan_core.processing.cpp_ingestion.models.ChunkGraph.add_edge", "auto_param_L173_36", 0.0)  
5824:  
5825:  
5826: @dataclass  
5827: class IntegrityIndex:  
5828:     """  
5829:         Cryptographic integrity verification data.  
5830:  
5831:         Uses BLAKE2b (not BLAKE3) for aggregate hash computation.  
5832:         Implementation uses hashlib.blake2b over JSON-serialized chunk hashes.  
5833:     """  
5834:     blake2b_root: str # Aggregate hash (BLAKE2b-256) of all chunk hashes  
5835:     chunk_hashes: dict[str, str] = field(default_factory=dict)  
5836:  
5837:  
5838: @dataclass  
5839: class CanonPolicyPackage:  
5840:     """  
5841:         Canon Policy Package - Complete output from phase-one ingestion.  
5842:  
5843:         This is the top-level container for all ingestion results.  
5844:         Also known as Smart Policy Chunks (SPC) in newer terminology.  
5845:     """  
5846:     schema_version: str  
5847:     chunk_graph: ChunkGraph  
5848:  
5849:     # Optional high-level metadata  
5850:     policy_manifest: PolicyManifest | None = None  
5851:     quality_metrics: QualityMetrics | None = None
```

```
5852:     integrity_index: IntegrityIndex | None = None
5853:
5854:     # Raw metadata
5855:     metadata: dict[str, Any] = field(default_factory=dict)
5856:
5857:
5858:
5859: =====
5860: FILE: src/farfán_pipeline/processing/policy_processor.py
5861: =====
5862:
5863: """
5864: Causal Framework Policy Plan Processor - Industrial Grade
5865: =====
5866:
5867: A mathematically rigorous, production-hardened system for extracting and
5868: validating causal evidence from Colombian local development plans against
5869: the DECALOGO framework's six-dimensional evaluation criteria.
5870:
5871: Architecture:
5872:     - Bayesian evidence accumulation for probabilistic confidence scoring
5873:     - Multi-scale text segmentation with coherence-preserving boundaries
5874:     - Differential privacy-aware pattern matching for reproducibility
5875:     - Entropy-based relevance ranking with TF-IDF normalization
5876:     - Graph-theoretic dependency validation for causal chain integrity
5877:
5878: Version: 3.0.0 | ISO 9001:2015 Compliant
5879: Author: Policy Analytics Research Unit
5880: License: Proprietary
5881: """
5882:
5883: import logging
5884: import re
5885: import unicodedata
5886: from collections import defaultdict
5887: from dataclasses import asdict, dataclass, field
5888: from enum import Enum
5889: from functools import lru_cache
5890: from pathlib import Path
5891: from typing import Any, ClassVar, Optional
5892:
5893: import numpy as np
5894:
5895: # Import runtime error fixes for defensive programming
5896: from farfan_pipeline.utils.runtime_error_fixes import ensure_list_return
5897:
5898: from farfan_pipeline.analysis.financiero_viability_tablas import PDETAnalysisException, QualityScore
5899: from farfan_pipeline.core.parameters import ParameterLoaderV2
5900: from farfan_pipeline.core.calibration.decorators import calibrated_method
5901: from farfan_pipeline.core.ports import (
5902:     PortDocumentLoader,
5903:     PortMunicipalOntology,
5904:     PortSemanticAnalyzer,
5905:     PortPerformanceAnalyzer,
5906:     PortContradictionDetector,
5907:     PortTemporalLogicVerifier,
```

```
5908:     PortBayesianConfidenceCalculator,
5909:     PortMunicipalAnalyzer,
5910: )
5911:
5912: try:
5913:     from farfan_pipeline.analysis.contradiction_deteccion import (
5914:         PolicyDimension as ContradictionPolicyDimension,
5915:     )
5916:     CONTRADICTION_MODULE_AVAILABLE = True
5917: except Exception as import_error:
5918:     CONTRADICTION_MODULE_AVAILABLE = False
5919:     logger = logging.getLogger(__name__)
5920:     logger.warning(
5921:         "Falling back to lightweight contradiction components due to import error: %s",
5922:         import_error,
5923:     )
5924:
5925: class ContradictionPolicyDimension(Enum): # type: ignore[misc]
5926:     DIAGNOSTICO = "diagnÃ³stico"
5927:     ESTRATEGICO = "estratÃ©gico"
5928:     PROGRAMATICO = "programÃ¡tico"
5929:     FINANCIERO = "plan plurianual de inversiones"
5930:     SEGUIMIENTO = "seguimiento y evaluaciÃ³n"
5931:     TERRITORIAL = "ordenamiento territorial"
5932:
5933:
5934: class _FallbackBayesianCalculator:
5935:     """Fallback Bayesian calculator when advanced module is unavailable."""
5936:
5937:     def __init__(self) -> None:
5938:         self.prior_alpha = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackBayesianCalculator.__init__", "auto_param_L64_31", 1.0)
5939:         self.prior_beta = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackBayesianCalculator.__init__", "auto_param_L65_30", 1.0)
5940:
5941:     def calculate_posterior(
5942:         self, evidence_strength: float, observations: int, domain_weight: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackBa-
5943:             yesianCalculator.__init__", "auto_param_L68_86", 1.0)
5944:     ) -> float:
5945:         alpha_post = self.prior_alpha + evidence_strength * observations * domain_weight
5946:         beta_post = self.prior_beta + (1 - evidence_strength) * observations * domain_weight
5947:         return alpha_post / (alpha_post + beta_post)
5948:
5949: class _FallbackTemporalVerifier:
5950:     """Fallback temporal verifier providing graceful degradation."""
5951:
5952:     @calibrated_method("farfan_core.processing.policy_processor._FallbackTemporalVerifier.verify_temporal_consistency")
5953:     def verify_temporal_consistency(self, statements: list[Any]) -> tuple[bool, list[dict[str, Any]]]:
5954:         return True, []
5955:
5956:
5957: class _FallbackContradictionDetector:
5958:     """Fallback contradiction detector providing graceful degradation."""
5959:
5960:     def detect(
5961:         self,
5962:         text: str,
```



```

6019:     "recursos_programaticos": [
6020:         r"\b(?:presupuesto\s+(?:plurianual|de\s+inversi[Ã³n])\b",
6021:         r"\b(?:plan\s+(?:plurianual|financiero|operativo\s+anual))\b",
6022:         r"\b(?:marco\s+fiscal\s+de\s+mediano\s+plazo)\b",
6023:         r"\b(?:trazabilidad\s+(?:presupuestal|program[Ã;a]tica))\b",
6024:     ],
6025:     "capacidad_institucional": [
6026:         r"\b(?:capacidad(?:es)?\s+(?:institucional(?:es)?|t[Ã©cnica(?:s)?))\b",
6027:         r"\b(?:talento\s+humano\s+(?:disponible|requerido))\b",
6028:         r"\b(?:gobernanza\s+(?:de\s+)?(?:datos|informaci[Ã³n])\b",
6029:         r"\b(?:brechas?\s+(?:de\s+)?implementaci[Ã³n])\b",
6030:     ],
6031: },
6032: CausalDimension.D2_ACTIVIDADES: {
6033:     "formalizacion_actividades": [
6034:         r"\b(?:plan\s+de\s+acci[Ã³n]\s+detallado)\b",
6035:         r"\b(?:matriz\s+de\s+(?:actividades|intervenciones))\b",
6036:         r"\b(?:cronograma\s+(?:de\s+)?ejecuci[Ã³n])\b",
6037:         r"\b(?:responsables?\s+(?:designados?|identificados?))\b",
6038:     ],
6039:     "mecanismo_causal": [
6040:         r"\b(?:mecanismo(?:s)?\s+causal(?:es))\b",
6041:         r"\b(?:teor[Ã-i]a\s+(?:de\s+)?intervenci[Ã³n])\b",
6042:         r"\b(?:cadena\s+(?:de\s+)?causaci[Ã³n])\b",
6043:         r"\b(?:v[Ã-i]nculo(?:s)?\s+explicativo(?:s))\b",
6044:     ],
6045:     "poblacion_objetivo": [
6046:         r"\b(?:poblaci[Ã³n]\s+(?:diana|objetivo|beneficiaria))\b",
6047:         r"\b(?:criterios?\s+de\s+focalizaci[Ã³n])\b",
6048:         r"\b(?:segmentaci[Ã³n]\s+(?:territorial|poblacional))\b",
6049:     ],
6050:     "dosificacion_intervencion": [
6051:         r"\b(?:dosificacion[Ã³n]\s+(?:de\s+)?(?:la\s+)?intervenci[Ã³n])\b",
6052:         r"\b(?:intensidad\s+(?:de\s+)?tratamiento)\b",
6053:         r"\b(?:duraci[Ã³n]\s+(?:de\s+)?exposici[Ã³n])\b",
6054:     ],
6055: },
6056: CausalDimension.D3_PRODUCTOS: {
6057:     "indicadores_producto": [
6058:         r"\b(?:indicador(?:es)?\s+de\s+(?:producto|output|gesti[Ã³n])\b",
6059:         r"\b(?:entregables?\s+verificables?)\b",
6060:         r"\b(?:metas?\s+(?:de\s+)?producto)\b",
6061:     ],
6062:     "verificabilidad": [
6063:         r"\b(?:f[Ã³rmula]\s+(?:de\s+)?(?:c[Ã;a]lculo|medici[Ã³n])\b",
6064:         r"\b(?:fuente(?:s)?\s+(?:de\s+)?verificaci[Ã³n])\b",
6065:         r"\b(?:medio(?:s)?\s+de\s+(?:prueba|evidencia))\b",
6066:     ],
6067:     "trazabilidad_producto": [
6068:         r"\b(?:trazabilidad\s+(?:de\s+)?productos?)\b",
6069:         r"\b(?:sistema\s+de\s+registro)\b",
6070:         r"\b(?:cobertura\s+(?:real|efectiva))\b",
6071:     ],
6072: },
6073: CausalDimension.D4_RESULTADOS: {
6074:     "metricas_outcome": [

```

```

6075:         r"\b(?:indicador(?:es)?|m[Ã©]trica(?:s)?)\s+de\s+(?:resultado|outcome))\b",
6076:         r"\b(?:criterios?\s+de\s+[Ã©]xito)\b",
6077:         r"\b(?:umbral(?:es)?\s+de\s+desempe[Ã±o])\b",
6078:     ],
6079:     "encadenamiento_causal": [
6080:         r"\b(?:encadenamiento\s+(?:causal|l[Ã³]gico))\b",
6081:         r"\b(?:ruta(?:s)?\s+cr[Ã-i]tica(?:s)?)\b",
6082:         r"\b(?:dependencias?\s+causales?)\b",
6083:     ],
6084:     "ventana_maduracion": [
6085:         r"\b(?:ventana\s+de\s+maduraci[Ã³]n)\b",
6086:         r"\b(?:horizonte\s+(?:de\s+)?resultados?)\b",
6087:         r"\b(?:rezago(?:s)?\s+(?:temporal(?:es)?|esperado(?:s)?))\b",
6088:     ],
6089:     "nivel_ambicion": [
6090:         r"\b(?:nivel\s+de\s+ambici[Ã³]n)\b",
6091:         r"\b(?:metas?\s+(?:incrementales?|transformacionales?))\b",
6092:     ],
6093: },
6094: CausalDimension.D5_IMPACTOS: {
6095:     "efectos_largo_plazo": [
6096:         r"\b(?:impacto(?:s)?\s+(:esperado(?:s)?|de\s+largo\s+plazo))\b",
6097:         r"\b(?:efectos?\s+(?:sostenidos?|duraderos?))\b",
6098:         r"\b(?:transformaci[Ã³]n\s+(?:estructural|sistÃ©mica))\b",
6099:     ],
6100:     "rutas_transmision": [
6101:         r"\b(?:ruta(?:s)?\s+de\s+transmisi[Ã³]n)\b",
6102:         r"\b(?:canales?\s+(?:de\s+)?(?:impacto|propagaci[Ã³]n))\b",
6103:         r"\b(?:efectos?\s+(?:directos?|indirectos?|multiplicadores?))\b",
6104:     ],
6105:     "proxies_mensurables": [
6106:         r"\b(?:proxies?\s+(?:de\s+)?impacto)\b",
6107:         r"\b(?:indicadores?\s+(?:compuestos?|s[Ã-i]ntesis))\b",
6108:         r"\b(?:medidas?\s+(?:indirectas?|aproximadas?))\b",
6109:     ],
6110:     "alineacion Marcos": [
6111:         r"\b(?:alineaci[Ã³]n\s+con\s+(?:PND|Plan\s+Nacional))\b",
6112:         r"\b(?:ODS\s+d+|Objetivo(?:s)?\s+de\s+Desarrollo\s+Sostenible)\b",
6113:         r"\b(?:coherencia\s+(?:vertical|horizontal))\b",
6114:     ],
6115: },
6116: CausalDimension.D6_CAUSALIDAD: {
6117:     "teoria_cambio_explicita": [
6118:         r"\b(?:teor[Ã-i]a\s+de(?:l)?\s+cambio)\b",
6119:         r"\b(?:modelo\s+l[Ã³]gico\s+(?:integrado|completo))\b",
6120:         r"\b(?:marco\s+causal\s+(?:expl[Ã-i]cito|formalizado))\b",
6121:     ],
6122:     "diagrama_causal": [
6123:         r"\b(?:diagrama\s+(?:causal|DAG)\s+de\s+flujo)\b",
6124:         r"\b(?:representaci[Ã³]n\s+gr[Ã;a]fica\s+causal)\b",
6125:         r"\b(?:mapa\s+(?:de\s+)?relaciones?)\b",
6126:     ],
6127:     "supuestos_verificables": [
6128:         r"\b(?:supuestos?\s+(?:verificables?|cr[Ã-i]ticos?))\b",
6129:         r"\b(?:hip[Ã³]tesis\s+(?:causales?|comprobables?))\b",
6130:         r"\b(?:condiciones?\s+(?:necesarias?|suficientes?))\b",

```

```
6131:     ],
6132:     "mediadores_moderadores": [
6133:         r"\b(?:mediador(?:es)?|moderador(?:es))\b",
6134:         r"\b(?:variables?s+(?:intermedias?|mediadoras?|moderadoras))\b",
6135:     ],
6136:     "validacion_logica": [
6137:         r"\b(?:validaci[Ã³]n\s+(?:l[Ã³]gica|emp[Ã–]rica))\b",
6138:         r"\b(?:pruebas?\s+(?:de\s+)?consistencia)\b",
6139:         r"\b(?:auditor[Ã–]a\s+causal)\b",
6140:     ],
6141:     "sistema_seguimiento": [
6142:         r"\b(?:sistema\s+de\s+(?:seguimiento|monitoreo))\b",
6143:         r"\b(?:tablero\s+de\s+(?:control|indicadores))\b",
6144:         r"\b(?:evaluaci[Ã³]n\s+(?:continua|peri[Ã³]dica))\b",
6145:     ],
6146: },
6147: },
6148:
6149: # =====
6150: # CONFIGURATION ARCHITECTURE
6151: # =====
6152:
6153: @dataclass(frozen=True)
6154: class ProcessorConfig:
6155:     """Immutable configuration for policy plan processing."""
6156:
6157:     preserve_document_structure: bool = True
6158:     enable_semantic_tagging: bool = True
6159:     confidence_threshold: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements",
"auto_param_L301_34", 0.65)
6160:     context_window_chars: int = 400
6161:     max_evidence_per_pattern: int = 5
6162:     enable_bayesian_scoring: bool = True
6163:     utf8_normalization_form: str = "NFC"
6164:
6165:     # Advanced controls
6166:     entropy_weight: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto
_param_L308_28", 0.3)
6167:     proximity_decay_rate: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements",
"auto_param_L309_34", 0.15)
6168:     min_sentence_length: int = 20
6169:     max_sentence_length: int = 500
6170:     bayesian_prior_confidence: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_stateme
nts", "auto_param_L312_39", 0.5)
6171:     bayesian_entropy_weight: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statem
ents", "auto_param_L313_37", 0.3)
6172:     minimum_dimension_scores: dict[str, float] = field(
6173:         default_factory=lambda:
6174:             "D1": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L31
6_18", 0.50),
6175:             "D2": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L31
7_18", 0.50),
6176:             "D3": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L31
8_18", 0.50),
6177:             "D4": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L31
9_18", 0.50),
```

```

6178:         "D5": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L32
0_18", 0.50),
6179:         "D6": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements", "auto_param_L32
1_18", 0.50),
6180:     }
6181:   )
6182:   critical_dimension_overrides: dict[str, float] = field(
6183:     default_factory=lambda: {"D1": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_stateme
nts", "auto_param_L325_39", 0.55), "D6": ParameterLoaderV2.get("farfan_core.processing.policy_processor._FallbackContradictionDetector._extract_policy_statements",
"auto_param_L325_51", 0.55)}
6184:   )
6185:   differential_focus_indicators: tuple[str, ...] = (
6186:     "enfoque diferencial",
6187:     "enfoque de gÃ©nero",
6188:     "mujeres rurales",
6189:     "poblaciÃ³n vÃ¡ctima",
6190:     "firmantes del acuerdo",
6191:     "comunidades indÃ³genas",
6192:     "poblaciÃ³n LGBTIQ+",
6193:     "juventud rural",
6194:     "comunidades ribereÃ±as",
6195:   )
6196:   adaptability_indicators: tuple[str, ...] = (
6197:     "mecanismo de ajuste",
6198:     "retroalimentaciÃ³n",
6199:     "aprendizaje",
6200:     "monitoreo adaptativo",
6201:     "ciclo de mejora",
6202:     "sistema de alerta temprana",
6203:     "evaluaciÃ³n continua",
6204:   )
6205:
6206:   LEGACY_PARAM_MAP: ClassVar[dict[str, str]] = {
6207:     "keep_structure": "preserve_document_structure",
6208:     "tag_elements": "enable_semantic_tagging",
6209:     "threshold": "confidence_threshold",
6210:   }
6211:
6212:   @classmethod
6213:   def from_legacy(cls, **kwargs: Any) -> "ProcessorConfig":
6214:     """Construct configuration from legacy parameter names."""
6215:     normalized = {}
6216:     for key, value in kwargs.items():
6217:       canonical = cls.LEGACY_PARAM_MAP.get(key, key)
6218:       normalized[canonical] = value
6219:     return cls(**normalized)
6220:
6221:   @calibrated_method("farfan_core.processing.policy_processor.ProcessorConfig.validate")
6222:   def validate(self) -> None:
6223:     """Validate configuration parameters."""
6224:     if not ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L366_15", 0.0) <= self.confidence_thres
hold <= ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L366_51", 1.0):
6225:       raise ValueError("confidence_threshold must be in [0, 1]")
6226:     if self.context_window_chars < 100:
6227:       raise ValueError("context_window_chars must be >= 100")
6228:     if self.entropy_weight < 0 or self.entropy_weight > 1:

```

```

6229:             raise ValueError("entropy_weight must be in [0, 1]")
6230:             if not ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L372_15", 0.0) <= self.bayesian_prior_c
onfidence <= ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L372_56", 1.0):
6231:                 raise ValueError("bayesian_prior_confidence must be in [0, 1]")
6232:                 if not ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L374_15", 0.0) <= self.bayesian_entropy
_weight <= ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L374_54", 1.0):
6233:                     raise ValueError("bayesian_entropy_weight must be in [0, 1]")
6234:                     for dimension, threshold in self.minimum_dimension_scores.items():
6235:                         if not ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L377_19", 0.0) <= threshold <= Para
meterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L377_39", 1.0):
6236:                             raise ValueError(
6237:                                 f"minimum_dimension_scores[{dimension}] must be in [0, 1]"
6238:                             )
6239:                             for dimension, threshold in self.critical_dimension_overrides.items():
6240:                                 if not ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L382_19", 0.0) <= threshold <= Para
meterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L382_39", 1.0):
6241:                                     raise ValueError(
6242:                                         f"critical_dimension_overrides[{dimension}] must be in [0, 1]"
6243:                                         )
6244:
6245: # =====
6246: # MATHEMATICAL SCORING ENGINE
6247: # =====
6248:
6249: class BayesianEvidenceScorer:
6250:     """
6251:         Bayesian evidence accumulation with entropy-weighted confidence scoring.
6252:
6253:         Implements a modified Dempster-Shafer framework for multi-evidence fusion
6254:         with automatic calibration against ground-truth policy corpora.
6255:     """
6256:
6257:     def __init__(
6258:         self,
6259:         prior_confidence: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L401_34", 0.5),
6260:         entropy_weight: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L402_32", 0.3),
6261:         calibration: dict[str, Any] | None = None,
6262:     ) -> None:
6263:         self.prior = prior_confidence
6264:         self.entropy_weight = entropy_weight
6265:         self._evidence_cache: dict[str, float] = {}
6266:         self.calibration = calibration or {}
6267:
6268:         # Defaults that can be overridden by calibration manifests
6269:         self.epsilon_clip: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L411_35", 0.02)
6270:         self.duplicate_gamma: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L412_38", 1.0)
6271:         self.cross_type_floor: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L413_39", 0.0)
6272:         self.source_quality_weights: dict[str, float] = {}
6273:         self.sector_multipliers: dict[str, float] = {}
6274:         self.sector_default: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L416_37", 1.0)
6275:         self.municipio_multipliers: dict[str, float] = {}
6276:         self.municipio_default: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.ProcessorConfig.validate", "auto_param_L418_40", 1.0)
6277:
6278:         self._configure_from_calibration()
6279:
6280:     @calibrated_method("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration")

```

```

6281:     def _configure_from_calibration(self) -> None:
6282:         config = self.calibration.get("bayesian_inference_robust") if isinstance(self.calibration, dict) else {}
6283:         if not isinstance(config, dict):
6284:             return
6285:
6286:         evidence_cfg = config.get("mechanistic_evidence_system", {})
6287:         if isinstance(evidence_cfg, dict):
6288:             stability = evidence_cfg.get("stability_controls", {})
6289:             if isinstance(stability, dict):
6290:                 self.epsilon_clip = float(stability.get("epsilon_clip", self.epsilon_clip))
6291:                 self.duplicate_gamma = float(stability.get("duplicate_gamma", self.duplicate_gamma))
6292:                 self.cross_type_floor = float(stability.get("cross_type_floor", self.cross_type_floor))
6293:                 self.epsilon_clip = min(max(self.epsilon_clip, ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L435_63", 0.0)), ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L435_69", 0.45))
6294:                 self.duplicate_gamma = max(ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L436_43", 0.0), self.duplicate_gamma)
6295:                 self.cross_type_floor = max(ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L437_44", 0.0), min(ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L437_53", 1.0), self.cross_type_floor))
6296:
6297:             weights = evidence_cfg.get("source_quality_weights", {})
6298:             if isinstance(weights, dict):
6299:                 self.source_quality_weights = {str(k): float(v) for k, v in weights.items() if isinstance(v, (int, float))}

6300:
6301:             context_cfg = config.get("theoretically_grounded_priors", {})
6302:             if isinstance(context_cfg, dict):
6303:                 hierarchy = context_cfg.get("hierarchical_context_priors", {})
6304:                 if isinstance(hierarchy, dict):
6305:                     sector = hierarchy.get("sector_multipliers", {})
6306:                     if isinstance(sector, dict):
6307:                         self.sector_multipliers = {str(k).lower(): float(v) for k, v in sector.items() if isinstance(v, (int, float))}
6308:                         self.sector_default = float(self.sector_multipliers.get("default", ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L450_87", 1.0)))
6309:                         muni = hierarchy.get("municipio_tamano_multipliers", {})
6310:                         if isinstance(muni, dict):
6311:                             self.municipio_multipliers = {str(k).lower(): float(v) for k, v in muni.items() if isinstance(v, (int, float))}
6312:                             self.municipio_default = float(self.municipio_multipliers.get("default", ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L454_93", 1.0)))
6313:
6314:             def compute_evidence_score(
6315:                 self,
6316:                 matches: list[str],
6317:                 total_corpus_size: int,
6318:                 pattern_specificity: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L460_37", 0.8),
6319:                 **kwargs: Any
6320:             ) -> float:
6321:                 """
6322:                     Compute probabilistic confidence score for evidence matches.
6323:
6324:                 Args:
6325:                     matches: List of matched text segments
6326:                     total_corpus_size: Total document size in characters
6327:                     pattern_specificity: Pattern discrimination power [0,1]
6328:                     **kwargs: Additional optional parameters for compatibility

```

```

6329:
6330:     Returns:
6331:         Calibrated confidence score in [0, 1]
6332:     """
6333:     if not matches:
6334:         return ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L476_19",
0.0)
6335:
6336:     # Term frequency normalization
6337:     tf = len(matches) / max(1, total_corpus_size / 1000)
6338:     if self.cross_type_floor:
6339:         tf = max(self.cross_type_floor, tf)
6340:
6341:     # Entropy-based diversity penalty
6342:     match_lengths = np.array([len(m) for m in matches])
6343:     entropy = self._calculate_shannon_entropy(match_lengths)
6344:
6345:     # Bayesian update
6346:     clip_low = self.epsilon_clip
6347:     clip_high = ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L489_20"
, 1.0) - self.epsilon_clip
6348:     pattern_specificity = max(clip_low, min(clip_high, pattern_specificity))
6349:
6350:     likelihood = min(ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L49
2_25", 1.0), tf * pattern_specificity)
6351:     posterior = (likelihood * self.prior) /
6352:         (likelihood * self.prior) + ((1 - likelihood) * (1 - self.prior))
6353:     )
6354:
6355:     # Entropy-weighted adjustment
6356:     final_score = (1 - self.entropy_weight) * posterior + self.entropy_weight * (
6357:         1 - entropy
6358:     )
6359:
6360:     # Apply duplicate penalty if provided by caller
6361:     if kwargs.get("duplicate_penalty"):
6362:         final_score *= self.duplicate_gamma
6363:
6364:     # Apply source quality weighting
6365:     if self.source_quality_weights:
6366:         source_quality = kwargs.get("source_quality")
6367:         if source_quality is not None:
6368:             weight = self._lookup_weight(self.source_quality_weights, source_quality, default=ParameterLoaderV2.get("farfan_core.processing.policy_proce
ssor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L510_98", 1.0))
6369:             final_score *= weight
6370:
6371:     # Context multipliers (sector / municipality)
6372:     sector = kwargs.get("sector") or kwargs.get("policy_sector")
6373:     if self.sector_multipliers:
6374:         final_score *= self._lookup_weight(self.sector_multipliers, sector, default=self.sector_default)
6375:
6376:     municipio = kwargs.get("municipio_tamano") or kwargs.get("municipio_size")
6377:     if self.municipio_multipliers:
6378:         final_score *= self._lookup_weight(self.municipio_multipliers, municipio, default=self.municipio_default)
6379:
6380:     return np.clip(final_score, ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "aut

```

```
o_param_L522_36", 0.0), ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L522_41", 1
.0))
6381:
6382:     @staticmethod
6383:     def _calculate_shannon_entropy(values: np.ndarray, **kwargs: Any) -> float:
6384:         """Calculate normalized Shannon entropy for value distribution.
6385:
6386:         Args:
6387:             values: Array of numerical values
6388:             **kwargs: Additional optional parameters for compatibility
6389:
6390:         Returns:
6391:             Normalized Shannon entropy
6392:         """
6393:         if len(values) < 2:
6394:             return ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configure_from_calibration", "auto_param_L536_19",
0.0)
6395:
6396:         # Discrete probability distribution
6397:         hist, _ = np.histogram(values, bins=min(10, len(values)))
6398:         prob = hist / hist.sum()
6399:         prob = prob[prob > 0] # Remove zeros
6400:
6401:         entropy = -np.sum(prob * np.log2(prob))
6402:         max_entropy = np.log2(len(prob)) if len(prob) > 1 else ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._config
ure_from_calibration", "auto_param_L544_63", 1.0)
6403:
6404:         return entropy / max_entropy if max_entropy > 0 else ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidenceScorer._configur
e_from_calibration", "auto_param_L546_61", 0.0)
6405:
6406:     @staticmethod
6407:     def _lookup_weight(mapping: dict[str, float], key: Any, default: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.BayesianEvidence
Scorer._configure_from_calibration", "auto_param_L549_77", 1.0)) -> float:
6408:         if not mapping:
6409:             return default
6410:         if key is None:
6411:             return mapping.get("default", default)
6412:         if isinstance(key, str):
6413:             direct = mapping.get(key)
6414:             if direct is not None:
6415:                 return direct
6416:             lowered = key.lower()
6417:             for candidate, value in mapping.items():
6418:                 if isinstance(candidate, str) and candidate.lower() == lowered:
6419:                     return value
6420:         return mapping.get("default", default)
6421:
6422: # =====
6423: # ADVANCED TEXT PROCESSOR
6424: # =====
6425:
6426: class PolicyTextProcessor:
6427:     """
6428:         Industrial-grade text processing with multi-scale segmentation and
6429:         coherence-preserving normalization for policy document analysis.
6430:     """

```

```
6431:     def __init__(self, config: ProcessorConfig, *, calibration: dict[str, Any] | None = None) -> None:
6432:         self.config = config
6433:         self.calibration = calibration or {}
6434:         self._compiled_patterns: dict[str, re.Pattern] = {}
6435:         self._sentence_boundaries = re.compile(
6436:             r"(?<=[.!?])\s+(?=([A-ZÀ\201À\211À\215À\223À\232À\221]))|(?<=\n\n)"
6437:         )
6438:
6439:
6440:     @calibrated_method("farfan_core.processing.policy_processor.PolicyTextProcessor.normalize_unicode")
6441:     def normalize_unicode(self, text: str) -> str:
6442:         """Apply canonical Unicode normalization (NFC/NFKC)."""
6443:         return unicodedata.normalize(self.config.utf8_normalization_form, text)
6444:
6445:     @calibrated_method("farfan_core.processing.policy_processor.PolicyTextProcessor.segment_into_sentences")
6446:     def segment_into_sentences(self, text: str, **kwargs: Any) -> list[str]:
6447:         """
6448:             Segment text into sentences with context-aware boundary detection.
6449:             Handles abbreviations, numerical lists, and Colombian naming conventions.
6450:
6451:             Args:
6452:                 text: Input text to segment
6453:                 **kwargs: Additional optional parameters for compatibility
6454:
6455:             Returns:
6456:                 List of sentence strings
6457:         """
6458:         # Protect common abbreviations
6459:         protected = text
6460:         protected = re.sub(r"\bDr\b.", "Dr___", protected)
6461:         protected = re.sub(r"\bSr\b.", "Sr___", protected)
6462:         protected = re.sub(r"\bart\b.", "art___", protected)
6463:         protected = re.sub(r"\bInc\b.", "Inc___", protected)
6464:
6465:         sentences = self._sentence_boundaries.split(protected)
6466:
6467:         # Restore protected patterns
6468:         sentences = [s.replace("___", ".") for s in sentences]
6469:
6470:         # Filter by length constraints
6471:         return [
6472:             s.strip()
6473:             for s in sentences
6474:             if self.config.min_sentence_length
6475:                 <= len(s.strip())
6476:                 <= self.config.max_sentence_length
6477:         ]
6478:
6479:     def extract_contextual_window(
6480:         self, text: str, match_position: int, window_size: int
6481:     ) -> str:
6482:         """Extract semantically coherent context window around a match."""
6483:         start = max(0, match_position - window_size // 2)
6484:         end = min(len(text), match_position + window_size // 2)
6485:
6486:         # Expand to sentence boundaries
```

```
6487:         while start > 0 and text[start] not in ".!?\n":
6488:             start -= 1
6489:         while end < len(text) and text[end] not in ".!?\n":
6490:             end += 1
6491:
6492:         return text[start:end].strip()
6493:
6494:     @lru_cache(maxsize=256)
6495:     @calibrated_method("farfan_core.processing.policy_processor.PolicyTextProcessor.compile_pattern")
6496:     def compile_pattern(self, pattern_str: str) -> re.Pattern:
6497:         """Cache and compile regex patterns for performance."""
6498:         return re.compile(pattern_str, re.IGNORECASE | re.UNICODE)
6499:
6500: # =====
6501: # CORE INDUSTRIAL PROCESSOR
6502: # =====
6503:
6504: @dataclass
6505: class EvidenceBundle:
6506:     """Structured evidence container with provenance and confidence metadata."""
6507:
6508:     dimension: CausalDimension
6509:     category: str
6510:     matches: list[str] = field(default_factory=list)
6511:     confidence: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.PolicyTextProcessor.compile_pattern", "auto_param_L653_24", 0.0)
6512:     context_windows: list[str] = field(default_factory=list)
6513:     match_positions: list[int] = field(default_factory=list)
6514:
6515:     @calibrated_method("farfan_core.processing.policy_processor.EvidenceBundle.to_dict")
6516:     def to_dict(self) -> dict[str, Any]:
6517:         return {
6518:             "dimension": self.dimension.value,
6519:             "category": self.category,
6520:             "match_count": len(self.matches),
6521:             "confidence": round(self.confidence, 4),
6522:             "evidence_samples": self.matches[:3],
6523:             "context_preview": self.context_windows[:2],
6524:         }
6525:
6526: class IndustrialPolicyProcessor:
6527:     """
6528:     State-of-the-art policy plan processor implementing rigorous causal
6529:     framework analysis with Bayesian evidence scoring and graph-theoretic
6530:     validation for Colombian local development plans.
6531:
6532:     This processor provides core analysis capabilities for policy documents.
6533:
6534:     DEPRECATION NOTE: The questionnaire_path parameter is deprecated.
6535:     Modern pipelines use SPC (Smart Policy Chunks) ingestion which handles
6536:     questionnaire integration separately.
6537:     """
6538:
6539:     def __init__(
6540:         self,
6541:         config: ProcessorConfig | None = None,
6542:         questionnaire_path: Path | None = None, # DEPRECATED: Kept for API compatibility only
```

```
6543:     *,
6544:     ontology: PortMunicipalOntology | None = None,
6545:     semantic_analyzer: PortSemanticAnalyzer | None = None,
6546:     performance_analyzer: PortPerformanceAnalyzer | None = None,
6547:     contradiction_detector: PortContradictionDetector | None = None,
6548:     temporal_verifier: PortTemporalLogicVerifier | None = None,
6549:     confidence_calculator: PortBayesianConfidenceCalculator | None = None,
6550:     municipal_analyzer: PortMunicipalAnalyzer | None = None,
6551: ) -> None:
6552:     # DEPRECATION WARNING: questionnaire_path parameter is deprecated
6553:     if questionnaire_path is not None:
6554:         import warnings
6555:         warnings.warn(
6556:             "The 'questionnaire_path' parameter is deprecated and will be ignored. "
6557:             "Modern SPC pipelines handle questionnaire integration separately. "
6558:             "Use CPPIngestionPipeline instead.",
6559:             DeprecationWarning,
6560:             stacklevel=2
6561:         )
6562:
6563:     self.config = config or ProcessorConfig()
6564:     self.config.validate()
6565:
6566:     self.text_processor = PolicyTextProcessor(self.config)
6567:     self.scorer = BayesianEvidenceScorer(
6568:         prior_confidence=self.config.bayesian_prior_confidence,
6569:         entropy_weight=self.config.bayesian_entropy_weight,
6570:     )
6571:
6572:     if ontology is None or semantic_analyzer is None or performance_analyzer is None:
6573:         from farfan_pipeline.core.wiring.analysis_factory import (
6574:             create_municipal_ontology,
6575:             create_semantic_analyzer,
6576:             create_performance_analyzer,
6577:         )
6578:         ontology = ontology or create_municipal_ontology()
6579:         semantic_analyzer = semantic_analyzer or create_semantic_analyzer(ontology)
6580:         performance_analyzer = performance_analyzer or create_performance_analyzer(ontology)
6581:
6582:     if contradiction_detector is None:
6583:         from farfan_pipeline.core.wiring.analysis_factory import create_contradiction_detector
6584:         contradiction_detector = create_contradiction_detector()
6585:
6586:     if temporal_verifier is None:
6587:         from farfan_pipeline.core.wiring.analysis_factory import create_temporal_logic_verifier
6588:         temporal_verifier = create_temporal_logic_verifier()
6589:
6590:     if confidence_calculator is None:
6591:         from farfan_pipeline.core.wiring.analysis_factory import create_bayesian_confidence_calculator
6592:         confidence_calculator = create_bayesian_confidence_calculator()
6593:
6594:     if municipal_analyzer is None:
6595:         from farfan_pipeline.core.wiring.analysis_factory import create_municipal_analyzer
6596:         municipal_analyzer = create_municipal_analyzer()
6597:
6598:     self.ontology = ontology
```

```

6599:         self.semantic_analyzer = semantic_analyzer
6600:         self.performance_analyzer = performance_analyzer
6601:         self.contradiction_detector = contradiction_detector
6602:         self.temporal_verifier = temporal_verifier
6603:         self.confidence_calculator = confidence_calculator
6604:         self.municipal_analyzer = municipal_analyzer
6605:
6606:     # LEGACY: Questionnaire loading removed - this component is deprecated
6607:     # Modern SPC pipeline handles questionnaire injection separately
6608:     self.questionnaire_file_path = None
6609:     self.questionnaire_data = {"questions": []} # Empty stub for backward compatibility
6610:
6611:     # Compile pattern taxonomy
6612:     self._pattern_registry = self._compile_pattern_registry()
6613:
6614:     # Policy point keyword extraction
6615:     self.point_patterns: dict[str, re.Pattern] = {}
6616:     self._build_point_patterns()
6617:
6618:     # Processing statistics
6619:     self.statistics: dict[str, Any] = defaultdict(int)
6620:
6621: @calibrated_method("farfan_core.processing.policy_processor.IndustrialPolicyProcessor._load_questionnaire")
6622: def _load_questionnaire(self) -> dict[str, Any]:
6623:     """
6624:     LEGACY: Questionnaire loading disabled.
6625:
6626:     This method is kept for backward compatibility but returns empty data.
6627:     Modern SPC pipeline handles questionnaire injection separately.
6628:     """
6629:     logger.warning(
6630:         "IndustrialPolicyProcessor._load_questionnaire called but questionnaire "
6631:         "loading is disabled. This is a legacy component. Use SPC ingestion instead."
6632:     )
6633:     return {"questions": []}
6634:
6635: @calibrated_method("farfan_core.processing.policy_processor.IndustrialPolicyProcessor._compile_pattern_registry")
6636: def _compile_pattern_registry(self) -> dict[CausalDimension, dict[str, list[re.Pattern]]]:
6637:     """Compile all causal patterns into efficient regex objects."""
6638:     registry = {}
6639:     for dimension, categories in CAUSAL_PATTERN_TAXONOMY.items():
6640:         registry[dimension] = {}
6641:         for category, patterns in categories.items():
6642:             registry[dimension][category] = [
6643:                 self.text_processor.compile_pattern(p) for p in patterns
6644:             ]
6645:     return registry
6646:
6647: @calibrated_method("farfan_core.processing.policy_processor.IndustrialPolicyProcessor._build_point_patterns")
6648: def _build_point_patterns(self) -> None:
6649:     """
6650:     LEGACY: Pattern building from questionnaire disabled.
6651:
6652:     This method is kept for backward compatibility but does nothing.
6653:     Modern SPC pipeline handles question-aware chunking separately.
6654:

```

```

6655:     questions = self.questionnaire_data.get("questions", [])
6656:
6657:     if not questions:
6658:         logger.info(
6659:             "No questionnaire questions available. "
6660:             "This is expected for legacy IndustrialPolicyProcessor. "
6661:             "Use SPC ingestion for question-aware analysis."
6662:         )
6663:     return
6664:
6665:     # Legacy path (should not be reached in modern pipeline)
6666:     point_keywords: dict[str, set[str]] = defaultdict(set)
6667:
6668:     for question in questions:
6669:         point_code = question.get("point_code")
6670:         if not point_code:
6671:             continue
6672:
6673:         # Extract title keywords
6674:         title = question.get("point_title", "").lower()
6675:         if title:
6676:             point_keywords[point_code].add(title)
6677:
6678:         # Extract hint keywords (cleaned)
6679:         for hint in question.get("hints", []):
6680:             cleaned = re.sub(r"\b\w+\b", "", hint).strip().lower()
6681:             if len(cleaned) > 3:
6682:                 point_keywords[point_code].add(cleaned)
6683:
6684:         # Compile into optimized regex patterns
6685:     for point_code, keywords in point_keywords.items():
6686:         # Sort by length (prioritize longer phrases)
6687:         sorted_kw = sorted(keywords, key=len, reverse=True)
6688:         pattern_str = "|".join(rf"\b{re.escape(kw)}\b" for kw in sorted_kw if kw)
6689:         self.point_patterns[point_code] = re.compile(pattern_str, re.IGNORECASE)
6690:
6691:     logger.info(f"Compiled patterns for {len(self.point_patterns)} policy points")
6692:
6693:     @calibrated_method("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process")
6694:     def process(self, raw_text: str, **kwargs: Any) -> dict[str, Any]:
6695:         """
6696:             Execute comprehensive policy plan analysis.
6697:
6698:             Args:
6699:                 raw_text: Sanitized policy document text
6700:                 **kwargs: Additional optional parameters (e.g., text, sentences, tables) for compatibility
6701:
6702:             Returns:
6703:                 Structured analysis results with evidence bundles and confidence scores
6704:         """
6705:         if not raw_text or len(raw_text) < 100:
6706:             logger.warning("Input text too short for analysis")
6707:             return self._empty_result()
6708:
6709:         # Normalize and segment
6710:         normalized = self.text_processor.normalize_unicode(raw_text)

```

```
6711:     sentences = self.text_processor.segment_into_sentences(normalized)
6712:
6713:     logger.info(f"Processing document: {len(normalized)} chars, {len(sentences)} sentences")
6714:
6715:     # Extract metadata
6716:     metadata = self._extract_metadata(normalized)
6717:
6718:     # Evidence extraction by policy point
6719:     point_evidence = {}
6720:     for point_code in sorted(self.point_patterns.keys()):
6721:         evidence = self._extract_point_evidence(
6722:             normalized, sentences, point_code
6723:         )
6724:         if evidence:
6725:             point_evidence[point_code] = evidence
6726:
6727:     # Global causal dimension analysis
6728:     dimension_analysis = self._analyze_causal_dimensions(normalized, sentences)
6729:
6730:     # Semantic diagnostics and performance evaluation
6731:     semantic_cube = self.semantic_analyzer.extract_semantic_cube(sentences)
6732:     performance_analysis = self.performance_analyzer.analyze_performance(
6733:         semantic_cube
6734:     )
6735:
6736:     try:
6737:         contradiction_bundle = self._run_contradiction_analysis(normalized, metadata)
6738:     except PDETAAnalysisException as exc:
6739:         logger.error("Contradiction analysis failed: %s", exc)
6740:         contradiction_bundle = {
6741:             "reports": {},
6742:             "temporal_assessments": {},
6743:             "bayesian_scores": {},
6744:             "critical_diagnosis": {
6745:                 "critical_links": {},
6746:                 "risk_assessment": {},
6747:                 "intervention_recommendations": {}
6748:             },
6749:         },
6750:
6751:         quality_score = self._calculate_quality_score(
6752:             dimension_analysis, contradiction_bundle, performance_analysis
6753:         )
6754:
6755:         summary = self.municipal_analyzer._generate_summary(
6756:             semantic_cube,
6757:             performance_analysis,
6758:             contradiction_bundle["critical_diagnosis"],
6759:         )
6760:
6761:     # Compile results
6762:     return {
6763:         "metadata": metadata,
6764:         "point_evidence": point_evidence,
6765:         "dimension_analysis": dimension_analysis,
6766:         "semantic_cube": semantic_cube,
```

```

6767:         "performance_analysis": performance_analysis,
6768:         "critical_diagnosis": contradiction_bundle["critical_diagnosis"],
6769:         "contradiction_reports": contradiction_bundle["reports"],
6770:         "temporal_consistency": contradiction_bundle["temporal_assessments"],
6771:         "bayesian_dimension_scores": contradiction_bundle["bayesian_scores"],
6772:         "quality_score": asdict(quality_score),
6773:         "summary": summary,
6774:         "document_statistics": {
6775:             "character_count": len(normalized),
6776:             "sentence_count": len(sentences),
6777:             "point_coverage": len(point_evidence),
6778:             "avg_confidence": self._compute_avg_confidence(dimension_analysis),
6779:         },
6780:         "processing_status": "complete",
6781:         "config_snapshot": {
6782:             "confidence_threshold": self.config.confidence_threshold,
6783:             "bayesian_enabled": self.config.enable_bayesian_scoring,
6784:         },
6785:     }
6786:
6787:     def _match_patterns_in_sentences(
6788:         self, compiled_patterns: list, relevant_sentences: list[str], **kwargs: Any
6789:     ) -> tuple[list[str], list[int]]:
6790:         """
6791:             Execute pattern matching across relevant sentences and collect matches with positions.
6792:
6793:             Args:
6794:                 compiled_patterns: List of compiled regex patterns to match
6795:                 relevant_sentences: Filtered sentences to search within
6796:                 **kwargs: Additional optional parameters for compatibility
6797:
6798:             Returns:
6799:                 Tuple of (matched_strings, match_positions)
6800:         """
6801:         matches = []
6802:         positions = []
6803:
6804:         for compiled_pattern in compiled_patterns:
6805:             for sentence in relevant_sentences:
6806:                 for match in compiled_pattern.finditer(sentence):
6807:                     matches.append(match.group(0))
6808:                     positions.append(match.start())
6809:
6810:         return matches, positions
6811:
6812:     def _compute_evidence_confidence(
6813:         self, matches: list[str], text_length: int, pattern_specificity: float, **kwargs: Any
6814:     ) -> float:
6815:         """
6816:             Calculate confidence score for evidence based on pattern matches and contextual factors.
6817:
6818:             Args:
6819:                 matches: List of matched pattern strings
6820:                 text_length: Total length of the document text
6821:                 pattern_specificity: Specificity coefficient for pattern weighting
6822:                 **kwargs: Additional optional parameters for compatibility

```

```
6823:  
6824:     Returns:  
6825:         Computed confidence score  
6826:     """  
6827:     confidence = self.scorer.compute_evidence_score(  
6828:         matches, text_length, pattern_specificity=pattern_specificity  
6829:     )  
6830:     return confidence  
6831:  
6832:     def _construct_evidence_bundle(  
6833:         self,  
6834:         dimension: CausalDimension,  
6835:         category: str,  
6836:         matches: list[str],  
6837:         positions: list[int],  
6838:         confidence: float,  
6839:         **kwargs: Any  
6840:     ) -> dict[str, Any]:  
6841:         """  
6842:             Assemble evidence bundle from matched patterns and computed confidence.  
6843:  
6844:             Args:  
6845:                 dimension: Causal dimension classification  
6846:                 category: Specific category within dimension  
6847:                 matches: List of matched pattern strings  
6848:                 positions: List of match positions in text  
6849:                 confidence: Computed confidence score  
6850:                 **kwargs: Additional optional parameters for compatibility  
6851:  
6852:             Returns:  
6853:                 Serialized evidence bundle dictionary  
6854:             """  
6855:             bundle = EvidenceBundle(  
6856:                 dimension=dimension,  
6857:                 category=category,  
6858:                 matches=matches[: self.config.max_evidence_per_pattern],  
6859:                 confidence=confidence,  
6860:                 match_positions=positions[: self.config.max_evidence_per_pattern],  
6861:             )  
6862:             return bundle.to_dict()  
6863:  
6864:     def _run_contradiction_analysis(  
6865:         self, text: str, metadata: dict[str, Any]  
6866:     ) -> dict[str, Any]:  
6867:         """Execute contradiction and temporal diagnostics across all dimensions."""  
6868:  
6869:         if not self.contradiction_detector:  
6870:             raise PDETAnalysisException("Contradiction detector unavailable")  
6871:  
6872:         plan_name = metadata.get("title", "Plan de Desarrollo")  
6873:         dimension_mapping = {  
6874:             CausalDimension.D1_INSUMOS: ContradictionPolicyDimension.DIAGNOSTICO,  
6875:             CausalDimension.D2_ACTIVIDADES: ContradictionPolicyDimension.ESTRATEGICO,  
6876:             CausalDimension.D3_PRODUCTOS: ContradictionPolicyDimension.PROGRAMATICO,  
6877:             CausalDimension.D4_RESULTADOS: ContradictionPolicyDimension.SEGUIMIENTO,  
6878:             CausalDimension.D5_IMPACTOS: ContradictionPolicyDimension.TERRITORIAL,
```

```

6879:             CausalDimension.D6_CAUSALIDAD: ContradictionPolicyDimension.ESTRATEGICO,
6880:         }
6881:
6882:         domain_weights = {
6883:             CausalDimension.D1_INSUMOS: 1.1,
6884:             CausalDimension.D2_ACTIVIDADES: ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L
1000_44", 1.0),
6885:             CausalDimension.D3_PRODUCTOS: ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L10
01_42", 1.0),
6886:             CausalDimension.D4_RESULTADOS: 1.1,
6887:             CausalDimension.D5_IMPACTOS: 1.15,
6888:             CausalDimension.D6_CAUSALIDAD: 1.2,
6889:         }
6890:
6891:         reports: dict[str, Any] = {}
6892:         temporal_assessments: dict[str, Any] = {}
6893:         bayesian_scores: dict[str, float] = {}
6894:         critical_links: dict[str, Any] = {}
6895:         risk_assessment: dict[str, Any] = {}
6896:         intervention_recommendations: dict[str, Any] = {}
6897:
6898:         for dimension in CausalDimension:
6899:             policy_dimension = dimension_mapping.get(dimension)
6900:             try:
6901:                 report = self.contradiction_detector.detect(
6902:                     text, plan_name=plan_name, dimension=policy_dimension
6903:                 )
6904:             except Exception as exc: # pragma: no cover - external deps
6905:                 raise PDETAAnalysisException(
6906:                     f"Contradiction detection failed for {dimension.name}: {exc}"
6907:                 ) from exc
6908:
6909:             reports[dimension.value] = report
6910:
6911:             try:
6912:                 statements = self.contradiction_detector._extract_policy_statements( # type: ignore[attr-defined]
6913:                     text, policy_dimension
6914:                 )
6915:             except Exception: # pragma: no cover - best effort if detector lacks method
6916:                 statements = []
6917:
6918:             is_consistent, conflicts = self.temporal_verifier.verify_temporal_consistency(
6919:                 statements
6920:             )
6921:             temporal_assessments[dimension.value] = {
6922:                 "is_consistent": is_consistent,
6923:                 "conflicts": conflicts,
6924:             }
6925:
6926:             coherence_metrics = report.get("coherence_metrics", {})
6927:             coherence_score = float(coherence_metrics.get("coherence_score", ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicy
Processor.process", "auto_param_L1043_77", 0.0)))
6928:             observations = max(1, len(statements))
6929:             posterior = self.confidence_calculator.calculate_posterior(
6930:                 evidence_strength=max(coherence_score, ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "a
uto_param_L1046_55", 0.01)),

```

```

6931:             observations=observations,
6932:             domain_weight=domain_weights.get(dimension, ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process",
6933:                                         "auto_param_L1048_60", 1.0)),
6934:             bayesian_scores[dimension.value] = float(posterior)
6935:
6936:             total_contradictions = int(report.get("total_contradictions", 0))
6937:             if total_contradictions:
6938:                 keywords = []
6939:                 # Defensive: ensure contradictions is a list
6940:                 contradictions_list = ensure_list_return(report.get("contradictions", []))
6941:                 for contradiction in contradictions_list:
6942:                     ctype = contradiction.get("contradiction_type")
6943:                     if ctype:
6944:                         keywords.append(ctype)
6945:
6946:             severity = 1 - coherence_score if coherence_score else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProces-
sor.process", "auto_param_L1062_71", 0.5)
6947:             critical_links[dimension.value] = {
6948:                 "criticality_score": round(min(ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_
param_L1064_51", 1.0), max(ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1064_60", 0.0), severity
)), 4),
6949:                 "text_analysis": {
6950:                     "sentiment": "negative" if coherence_score < ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcesso-
r.process", "auto_param_L1066_69", 0.5) else "neutral",
6951:                     "keywords": keywords,
6952:                     "word_count": len(text.split()),
6953:                 },
6954:             }
6955:             risk_assessment[dimension.value] = {
6956:                 "overall_risk": "high" if total_contradictions > 3 else "medium",
6957:                 "risk_factors": keywords,
6958:             }
6959:             intervention_recommendations[dimension.value] = report.get(
6960:                 "recommendations", [])
6961:             )
6962:
6963:             return {
6964:                 "reports": reports,
6965:                 "temporal_assessments": temporal_assessments,
6966:                 "bayesian_scores": bayesian_scores,
6967:                 "critical_diagnosis": {
6968:                     "critical_links": critical_links,
6969:                     "risk_assessment": risk_assessment,
6970:                     "intervention_recommendations": intervention_recommendations,
6971:                 },
6972:             }
6973:
6974:             def _calculate_quality_score(
6975:                 self,
6976:                 dimension_analysis: dict[str, Any],
6977:                 contradiction_bundle: dict[str, Any],
6978:                 performance_analysis: dict[str, Any],
6979:             ) -> QualityScore:
6980:                 """Aggregate key indicators into a structured QualityScore dataclass."""
6981:

```

```
6982:     bayesian_scores = contradiction_bundle.get("bayesian_scores", {})
6983:     bayesian_values = list(bayesian_scores.values())
6984:     overall_score = float(np.mean(bayesian_values)) if bayesian_values else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1100_80", 0.0)
6985:
6986:     def _dimension_confidence(key: CausalDimension) -> float:
6987:         return float(
6988:             dimension_analysis.get(key.value, {}).get("dimension_confidence", ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1104_82", 0.0)))
6989:     )
6990:
6991:     temporal_flags = contradiction_bundle.get("temporal_assessments", {})
6992:     temporal_values = [
6993:         ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1109_12", 1.0) if assessment.get("is_consistent", True) else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1109_62", 0.0)
6994:         for assessment in temporal_flags.values()
6995:     ]
6996:     temporal_consistency = (
6997:         float(np.mean(temporal_values)) if temporal_values else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1113_68", 1.0)
6998:     )
6999:
7000:     reports = contradiction_bundle.get("reports", {})
7001:     coherence_scores = [
7002:         float(report.get("coherence_metrics", {}).get("coherence_score", ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1118_77", 0.0)))
7003:             for report in reports.values()
7004:     ]
7005:     causal_coherence = float(np.mean(coherence_scores)) if coherence_scores else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1121_85", 0.0)
7006:
7007:     objective_alignment = float(
7008:         reports.get(
7009:             CausalDimension.D4_RESULTADOS.value,
7010:             {}),
7011:         )
7012:         .get("coherence_metrics", {})
7013:         .get("objective_alignment", ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1129_40", 0.0))
7014:     )
7015:
7016:     confidence_interval = (
7017:         float(min(bayesian_values)) if bayesian_values else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1133_64", 0.0),
7018:         float(max(bayesian_values)) if bayesian_values else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1134_64", 0.0),
7019:     )
7020:
7021:     evidence = {
7022:         "bayesian_scores": bayesian_scores,
7023:         "dimension_confidences": {
7024:             key: value.get("dimension_confidence", ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1140_55", 0.0))
7025:                 for key, value in dimension_analysis.items()
7026:             },
7027:         "performance_metrics": performance_analysis.get("value_chain_metrics", {}),
```

```
7028:         }
7029:
7030:     return QualityScore(
7031:         overall_score=overall_score,
7032:         financial_feasibility=_dimension_confidence(CausalDimension.D1_INSUMOS),
7033:         indicator_quality=_dimension_confidence(CausalDimension.D3_PRODUCTOS),
7034:         responsibility_clarity=_dimension_confidence(CausalDimension.D2_ACTIVIDADES),
7035:         temporal_consistency=temporal_consistency,
7036:         pdet_alignment=objective_alignment,
7037:         causal_coherence=causal_coherence,
7038:         confidence_interval=confidence_interval,
7039:         evidence=evidence,
7040:     )
7041:
7042:     def _extract_point_evidence(
7043:         self, text: str, sentences: list[str], point_code: str
7044:     ) -> dict[str, Any]:
7045:         """Extract evidence for a specific policy point across all dimensions."""
7046:         pattern = self.point_patterns.get(point_code)
7047:         if not pattern:
7048:             return {}
7049:
7050:         # Find relevant sentences
7051:         relevant_sentences = [s for s in sentences if pattern.search(s)]
7052:         if not relevant_sentences:
7053:             return {}
7054:
7055:         # Search for dimensional evidence within relevant context
7056:         evidence_by_dimension = {}
7057:         for dimension, categories in self._pattern_registry.items():
7058:             dimension_evidence = []
7059:
7060:             for category, compiled_patterns in categories.items():
7061:                 matches, positions = self._match_patterns_in_sentences(
7062:                     compiled_patterns, relevant_sentences
7063:                 )
7064:
7065:                 if matches:
7066:                     confidence = self._compute_evidence_confidence(
7067:                         matches, len(text), pattern_specificity=ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.pro
cess", "auto_param_L1183_64", 0.85)
7068:                     )
7069:
7070:                     if confidence >= self.config.confidence_threshold:
7071:                         evidence_dict = self._construct_evidence_bundle(
7072:                             dimension, category, matches, positions, confidence
7073:                         )
7074:                         dimension_evidence.append(evidence_dict)
7075:
7076:             if dimension_evidence:
7077:                 evidence_by_dimension[dimension.value] = dimension_evidence
7078:
7079:         return evidence_by_dimension
7080:
7081:     def _analyze_causal_dimensions(
7082:         self, text: str, sentences: list[str] | None = None
```

```
7083:     ) -> dict[str, Any]:
7084:         """
7085:             Perform global analysis of causal dimensions across entire document.
7086:
7087:             Args:
7088:                 text: Full document text
7089:                 sentences: Optional pre-segmented sentences. If not provided, will be
7090:                             automatically extracted from text using the text processor.
7091:
7092:             Returns:
7093:                 Dictionary containing dimension scores and confidence metrics
7094:
7095:             Note:
7096:                 This function requires 'sentences' for optimal performance. If not provided,
7097:                 sentences will be extracted from text automatically, which may impact performance.
7098:             """
7099:             # Defensive validation: ensure sentences parameter is provided
7100:             if sentences is None:
7101:                 logger.warning(
7102:                     "_analyze_causal_dimensions called without 'sentences' parameter. "
7103:                     "Automatically extracting sentences from text. "
7104:                     "Expected signature: _analyze_causal_dimensions(self, text: str, sentences: List[str])"
7105:                 )
7106:             # Auto-extract sentences if not provided
7107:             sentences = self.text_processor.segment_into_sentences(text)
7108:
7109:             dimension_scores = {}
7110:
7111:             for dimension, categories in self._pattern_registry.items():
7112:                 total_matches = 0
7113:                 category_results = {}
7114:
7115:                 for category, compiled_patterns in categories.items():
7116:                     matches = []
7117:                     for pattern in compiled_patterns:
7118:                         for sentence in sentences:
7119:                             matches.extend(pattern.findall(sentence))
7120:
7121:                     if matches:
7122:                         confidence = self.scorer.compute_evidence_score(
7123:                             matches, len(text), pattern_specificity=ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.pro
cess", "auto_param_L1239_64", 0.80)
7124:                         )
7125:                         category_results[category] = {
7126:                             "match_count": len(matches),
7127:                             "confidence": round(confidence, 4),
7128:                         }
7129:                         total_matches += len(matches)
7130:
7131:                         dimension_scores[dimension.value] = {
7132:                             "categories": category_results,
7133:                             "total_matches": total_matches,
7134:                             "dimension_confidence": round(
7135:                                 np.mean([c["confidence"] for c in category_results.values()])
7136:                                 if category_results
7137:                                 else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1253_25", 0.0),
```

```
7138:         4,
7139:     ),
7140: }
7141:
7142:     return dimension_scores
7143:
7144: @staticmethod
7145: def _extract_metadata(text: str) -> dict[str, Any]:
7146:     """Extract key metadata from policy document header."""
7147:     # Title extraction
7148:     title_match = re.search(
7149:         r"(?i)plan\s+(?:de\s+)?desarrollo\s+(?:municipal|departamental|local)?\s*[:\n-]?\s*([^\n]{10,150})",
7150:         text[:2000],
7151:     )
7152:     title = title_match.group(1).strip() if title_match else "Sin título identificado"
7153:
7154:     # Entity extraction
7155:     entity_match = re.search(
7156:         r"(?i)(?:municipio|alcald[í]a|gobernaci[ó]n|distrito)\s+(?:de\s+)?([A-ZÀ\201À\211À\215À\223À\232À\221][a-zA-À-Ã°À-Ã±À\s]+)",
7157:         text[:3000],
7158:     )
7159:     entity = entity_match.group(1).strip() if entity_match else "Entidad no especificada"
7160:
7161:     # Period extraction
7162:     period_match = re.search(r"(20\d{2})\s*[-à\200\223â\200\224]\s*(20\d{2})", text[:3000])
7163:     period = {
7164:         "start_year": int(period_match.group(1)) if period_match else None,
7165:         "end_year": int(period_match.group(2)) if period_match else None,
7166:     }
7167:
7168:     return {
7169:         "title": title,
7170:         "entity": entity,
7171:         "period": period,
7172:         "extraction_timestamp": "2025-10-13",
7173:     }
7174:
7175: @staticmethod
7176: def _compute_avg_confidence(dimension_analysis: dict[str, Any]) -> float:
7177:     """Calculate average confidence across all dimensions."""
7178:     confidences = [
7179:         dim_data["dimension_confidence"]
7180:         for dim_data in dimension_analysis.values()
7181:         if dim_data.get("dimension_confidence", 0) > 0
7182:     ]
7183:     return round(np.mean(confidences), 4) if confidences else ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor.process", "auto_param_L1299_66", 0.0)
7184:
7185:     @calibrated_method("farfan_core.processing.policy_processor.IndustrialPolicyProcessor._empty_result")
7186:     def _empty_result(self) -> dict[str, Any]:
7187:         """Return structure for failed/empty processing."""
7188:         return {
7189:             "metadata": {},
7190:             "point_evidence": {},
7191:             "dimension_analysis": {},
7192:             "document_statistics": {}
```

```
7193:             "character_count": 0,
7194:             "sentence_count": 0,
7195:             "point_coverage": 0,
7196:             "avg_confidence": ParameterLoaderV2.get("farfan_core.processing.policy_processor.IndustrialPolicyProcessor._empty_result", "auto_param_L1312
_34", 0.0),
7197:         },
7198:         "processing_status": "failed",
7199:         "error": "Insufficient input for analysis",
7200:     }
7201:
7202:     def export_results(
7203:         self, results: dict[str, Any], output_path: str | Path
7204:     ) -> None:
7205:         """Export analysis results to JSON with formatted output."""
7206:         # Delegate to factory for I/O operation
7207:         from farfan_pipeline.processing.factory import save_json
7208:
7209:         save_json(results, output_path)
7210:         logger.info(f"Results exported to {output_path}")
7211:
7212: # =====
7213: # ENHANCED SANITIZER WITH STRUCTURE PRESERVATION
7214: # =====
7215:
7216: class AdvancedTextSanitizer:
7217:     """
7218:         Sophisticated text sanitization preserving semantic structure and
7219:         critical policy elements with differential privacy guarantees.
7220:     """
7221:
7222:     def __init__(self, config: ProcessorConfig) -> None:
7223:         self.config = config
7224:         self.protection_markers: dict[str, tuple[str, str]] = {
7225:             "heading": ("__HEAD_START__", "__HEAD_END__"),
7226:             "list_item": ("__LIST_START__", "__LIST_END__"),
7227:             "table_cell": ("__TABLE_START__", "__TABLE_END__"),
7228:             "citation": ("__CITE_START__", "__CITE_END__"),
7229:         }
7230:
7231:         @calibrated_method("farfan_core.processing.policy_processor.AdvancedTextSanitizer.sanitize")
7232:         def sanitize(self, raw_text: str) -> str:
7233:             """
7234:                 Execute comprehensive text sanitization pipeline.
7235:
7236:                 Pipeline stages:
7237:                 1. Unicode normalization (NFC)
7238:                 2. Structure element protection
7239:                 3. Whitespace normalization
7240:                 4. Special character handling
7241:                 5. Encoding validation
7242:             """
7243:             if not raw_text:
7244:                 return ""
7245:
7246:             # Stage 1: Unicode normalization
7247:             text = unicodedata.normalize(self.config.utf8_normalization_form, raw_text)
```

```
7248:  
7249:     # Stage 2: Protect structural elements  
7250:     if self.config.preserve_document_structure:  
7251:         text = self._protect_structure(text)  
7252:  
7253:     # Stage 3: Whitespace normalization  
7254:     text = re.sub(r"\t+", " ", text)  
7255:     text = re.sub(r"\n{3,}", "\n\n", text)  
7256:  
7257:     # Stage 4: Remove control characters (except newlines/tabs)  
7258:     text = "".join(  
7259:         char for char in text  
7260:         if unicodedata.category(char)[0] != "C" or char in "\n\t"  
7261:     )  
7262:  
7263:     # Stage 5: Restore protected elements  
7264:     if self.config.preserve_document_structure:  
7265:         text = self._restore_structure(text)  
7266:  
7267:     return text.strip()  
7268:  
7269: @calibrated_method("farfan_core.processing.policy_processor.AdvancedTextSanitizer._protect_structure")  
7270: def _protect_structure(self, text: str) -> str:  
7271:     """Mark structural elements for protection during sanitization."""  
7272:     protected = text  
7273:  
7274:     # Protect headings (numbered or capitalized lines)  
7275:     heading_pattern = re.compile(  
7276:         r"^(?:[\d]+\s+)?([A-ZÀ\201À\211À\215À\223À\232À\221] [A-ZÀ\201À\211À\215À\223À\232À\221a-zÀ;À©À-À³ÀºÀ±\s]{5,80})$",  
7277:         re.MULTILINE,  
7278:     )  
7279:     for match in reversed(list(heading_pattern.finditer(protected))):  
7280:         start, end = match.span()  
7281:         heading_text = match.group(0)  
7282:         protected = (  
7283:             protected[:start]  
7284:             + f"{self.protection_markers['heading'][0]}{heading_text}{self.protection_markers['heading'][1]}"  
7285:             + protected[end:]  
7286:         )  
7287:  
7288:     # Protect list items  
7289:     list_pattern = re.compile(r"^\s*[â\200$-\*\d]+[.\.])\s+(.+)$", re.MULTILINE)  
7290:     for match in reversed(list(list_pattern.finditer(protected))):  
7291:         start, end = match.span()  
7292:         item_text = match.group(0)  
7293:         protected = (  
7294:             protected[:start]  
7295:             + f"{self.protection_markers['list_item'][0]}{item_text}{self.protection_markers['list_item'][1]}"  
7296:             + protected[end:]  
7297:         )  
7298:  
7299:     return protected  
7300:  
7301: @calibrated_method("farfan_core.processing.policy_processor.AdvancedTextSanitizer._restore_structure")  
7302: def _restore_structure(self, text: str) -> str:  
7303:     """Remove protection markers after sanitization."""
```

```
7304:         restored = text
7305:         for _marker_type, (start_mark, end_mark) in self.protection_markers.items():
7306:             restored = restored.replace(start_mark, "")
7307:             restored = restored.replace(end_mark, "")
7308:         return restored
7309:
7310: # =====
7311: # INTEGRATED FILE HANDLING WITH RESILIENCE
7312: # =====
7313:
7314: class ResilientFileHandler:
7315:     """
7316:         Production-grade file I/O with automatic encoding detection,
7317:         retry logic, and comprehensive error classification.
7318:     """
7319:
7320:     ENCODINGS = ["utf-8", "utf-8-sig", "latin-1", "cp1252", "iso-8859-1"]
7321:
7322:     @classmethod
7323:     def read_text(cls, file_path: str | Path) -> str:
7324:         """
7325:             Read text file with automatic encoding detection and fallback cascade.
7326:
7327:             Args:
7328:                 file_path: Path to input file
7329:
7330:             Returns:
7331:                 Decoded text content
7332:
7333:             Raises:
7334:                 IOError: If file cannot be read with any supported encoding
7335:         """
7336:         # Delegate to factory for I/O operation
7337:         from farfan_pipeline.processing.factory import read_text_file
7338:
7339:         try:
7340:             return read_text_file(file_path, encodings=list(cls.ENCODINGS))
7341:         except Exception as e:
7342:             raise OSError(f"Failed to read {file_path} with any supported encoding") from e
7343:
7344:     @classmethod
7345:     def write_text(cls, content: str, file_path: str | Path) -> None:
7346:         """Write text content with UTF-8 encoding and directory creation."""
7347:         # Delegate to factory for I/O operation
7348:         from farfan_pipeline.processing.factory import write_text_file
7349:
7350:         write_text_file(content, file_path)
7351:
7352: # =====
7353: # UNIFIED ORCHESTRATOR
7354: # =====
7355:
7356: class PolicyAnalysisPipeline:
7357:     """
7358:         End-to-end orchestrator for Colombian local development plan analysis
7359:         implementing the complete DECALOGO causal framework evaluation workflow.
```

```
7360:
7361:    DEPRECATION NOTE: The questionnaire_path parameter is deprecated.
7362:    Modern pipelines use SPC (Smart Policy Chunks) ingestion which handles
7363:    questionnaire integration separately.
7364:    """
7365:
7366:    def __init__(
7367:        self,
7368:        config: ProcessorConfig | None = None,
7369:        questionnaire_path: Path | None = None, # DEPRECATED: Kept for API compatibility only
7370:    ) -> None:
7371:        # DEPRECATION WARNING: questionnaire_path parameter is deprecated
7372:        if questionnaire_path is not None:
7373:            import warnings
7374:            warnings.warn(
7375:                "The 'questionnaire_path' parameter is deprecated and will be ignored. "
7376:                "Modern SPC pipelines handle questionnaire integration separately. "
7377:                "Use CPPIngestionPipeline instead.",
7378:                DeprecationWarning,
7379:                stacklevel=2
7380:            )
7381:
7382:        self.config = config or ProcessorConfig()
7383:        self.sanitizer = AdvancedTextSanitizer(self.config)
7384:
7385:        from farfan_pipeline.core.wiring.analysis_factory import create_analysis_components
7386:
7387:        components = create_analysis_components()
7388:        self.document_loader = components['document_loader']
7389:        self.ontology = components['ontology']
7390:        self.semantic_analyzer = components['semantic_analyzer']
7391:        self.performance_analyzer = components['performance_analyzer']
7392:        self.temporal_verifier = components['temporal_verifier']
7393:        self.confidence_calculator = components['confidence_calculator']
7394:        self.contradiction_detector = components['contradiction_detector']
7395:        self.municipal_analyzer = components['municipal_analyzer']
7396:
7397:        self.processor = IndustrialPolicyProcessor(
7398:            self.config,
7399:            questionnaire_path,
7400:            ontology=self.ontology,
7401:            semantic_analyzer=self.semantic_analyzer,
7402:            performance_analyzer=self.performance_analyzer,
7403:            contradiction_detector=self.contradiction_detector,
7404:            temporal_verifier=self.temporal_verifier,
7405:            confidence_calculator=self.confidence_calculator,
7406:            municipal_analyzer=self.municipal_analyzer,
7407:        )
7408:        self.file_handler = ResilientFileHandler()
7409:
7410:    def analyze_file(
7411:        self,
7412:        input_path: str | Path,
7413:        output_path: str | Path | None = None,
7414:    ) -> dict[str, Any]:
7415:        """
```

```
7416:     Execute complete analysis pipeline on a policy document file.
7417:
7418:     Args:
7419:         input_path: Path to input policy document (text format)
7420:         output_path: Optional path for JSON results export
7421:
7422:     Returns:
7423:         Complete analysis results dictionary
7424:         """
7425:     input_path = Path(input_path)
7426:     logger.info(f"Starting analysis of {input_path}")
7427:
7428:     # Stage 1: Load document
7429:     raw_text = ""
7430:     suffix = input_path.suffix.lower()
7431:     if suffix == ".pdf":
7432:         raw_text = self.document_loader.load_pdf(str(input_path))
7433:     elif suffix in {".docx", ".doc"}:
7434:         raw_text = self.document_loader.load_docx(str(input_path))
7435:
7436:     if not raw_text:
7437:         raw_text = self.file_handler.read_text(input_path)
7438:     logger.info(f"Loaded {len(raw_text)} characters from {input_path.name}")
7439:
7440:     # Stage 2: Sanitize
7441:     sanitized_text = self.sanitizer.sanitize(raw_text)
7442:     reduction_pct = 100 * (1 - len(sanitized_text) / max(1, len(raw_text)))
7443:     logger.info(f"Sanitization: {reduction_pct:.1f}% size reduction")
7444:
7445:     # Stage 3: Process
7446:     results = self.processor.process(sanitized_text)
7447:     results["pipeline_metadata"] = {
7448:         "input_file": str(input_path),
7449:         "raw_size": len(raw_text),
7450:         "sanitized_size": len(sanitized_text),
7451:         "reduction_percentage": round(reduction_pct, 2),
7452:     }
7453:
7454:     # Stage 4: Export if requested
7455:     if output_path:
7456:         self.processor.export_results(results, output_path)
7457:
7458:     logger.info(f"Analysis complete: {results['processing_status']}")
7459:     return results
7460:
7461: @calibrated_method("farfan_core.processing.policy_processor.PolicyAnalysisPipeline.analyze_text")
7462: def analyze_text(self, raw_text: str) -> dict[str, Any]:
7463:     """
7464:     Execute analysis pipeline on raw text input.
7465:
7466:     Args:
7467:         raw_text: Raw policy document text
7468:
7469:     Returns:
7470:         Complete analysis results dictionary
7471:         """
```

```
7472:     sanitized_text = self.sanitizer.sanitize(raw_text)
7473:     return self.processor.process(sanitized_text)
7474:
7475: # =====
7476: # FACTORY FUNCTIONS FOR BACKWARD COMPATIBILITY
7477: # =====
7478:
7479: def create_policy_processor(
7480:     preserve_structure: bool = True,
7481:     enable_semantic_tagging: bool = True,
7482:     confidence_threshold: float = ParameterLoaderV2.get("farfan_core.processing.policy_processor.PolicyAnalysisPipeline.analyze_text", "auto_param_L1595_34"
7483: , 0.65),
7484:     **kwargs: Any,
7485: ) -> PolicyAnalysisPipeline:
7486: """
7487:     Factory function for creating policy analysis pipeline with legacy support.
7488:
7489:     Args:
7490:         preserve_structure: Enable document structure preservation
7491:         enable_semantic_tagging: Enable semantic element tagging
7492:         confidence_threshold: Minimum confidence threshold for evidence
7493:         **kwargs: Additional configuration parameters
7494:
7495:     Returns:
7496:         Configured PolicyAnalysisPipeline instance
7497: """
7498: config = ProcessorConfig(
7499:     preserve_document_structure=preserve_structure,
7500:     enable_semantic_tagging=enable_semantic_tagging,
7501:     confidence_threshold=confidence_threshold,
7502:     **kwargs,
7503: )
7504: return PolicyAnalysisPipeline(config=config)
7505: # =====
7506: # COMMAND-LINE INTERFACE
7507: # =====
7508:
7509: def main() -> None:
7510:     """Command-line interface for policy plan analysis."""
7511:     import argparse
7512:
7513:     parser = argparse.ArgumentParser(
7514:         description="Industrial-Grade Policy Plan Processor for Colombian Local Development Plans"
7515:     )
7516:     parser.add_argument("input_file", type=str, help="Input policy document path")
7517:     parser.add_argument(
7518:         "-o", "--output", type=str, help="Output JSON file path", default=None
7519:     )
7520:     parser.add_argument(
7521:         "-t",
7522:         "--threshold",
7523:         type=float,
7524:         default=ParameterLoaderV2.get("farfan_core.processing.policy_processor.PolicyAnalysisPipeline.analyze_text", "auto_param_L1637_16", 0.65),
7525:         help="Confidence threshold (0-1)",
7526:     )
```

```
7527:     parser.add_argument(
7528:         "-q",
7529:         "--questionnaire",
7530:         type=str,
7531:         help="Custom questionnaire JSON path",
7532:         default=None,
7533:     )
7534:     parser.add_argument(
7535:         "-v", "--verbose", action="store_true", help="Enable verbose logging"
7536:     )
7537:
7538:     args = parser.parse_args()
7539:
7540:     if args.verbose:
7541:         logging.getLogger().setLevel(logging.DEBUG)
7542:
7543:     # Configure and execute pipeline
7544:     config = ProcessorConfig(confidence_threshold=args.threshold)
7545:     questionnaire_path = Path(args.questionnaire) if args.questionnaire else None
7546:
7547:     pipeline = PolicyAnalysisPipeline(
7548:         config=config, questionnaire_path=questionnaire_path
7549:     )
7550:
7551:     try:
7552:         results = pipeline.analyze_file(args.input_file, args.output)
7553:
7554:         # Print summary
7555:         print("\n" + "=" * 70)
7556:         print("POLICY ANALYSIS SUMMARY")
7557:         print("=" * 70)
7558:         print(f"Document: {results['metadata'].get('title', 'N/A')}")
7559:         print(f"Entity: {results['metadata'].get('entity', 'N/A')}")
7560:         print(f"Period: {results['metadata'].get('period', {})}")
7561:         print(f"\nPolicy Points Covered: {results['document_statistics']['point_coverage']}")
7562:         print(f"Average Confidence: {results['document_statistics']['avg_confidence']:.2%}")
7563:         print(f"Total Sentences: {results['document_statistics']['sentence_count']}")
7564:         print("=" * 70 + "\n")
7565:
7566:     except Exception as e:
7567:         logger.error(f"Analysis failed: {e}", exc_info=True)
7568:         raise
7569:
7570:
7571:
7572: =====
7573: FILE: src/farfan_pipeline/processing/quality_gates.py
7574: =====
7575:
7576: """
7577: Quality gates for SPC (Smart Policy Chunks) validation - Canonical Phase-One.
7578:
7579: Validates quality metrics, enforces invariants, and ensures compatibility
7580: with downstream phases in the canonical pipeline flux.
7581:
7582: MAXIMUM STANDARD: No tolerance for data quality degradation.
```

```
7583: """
7584:
7585: import logging
7586: from pathlib import Path
7587: from typing import Any
7588: # from farfan_core import get_parameter_loader # CALIBRATION DISABLED
7589: from farfan_pipeline.core.calibration.decorators import calibrated_method
7590:
7591: logger = logging.getLogger(__name__)
7592:
7593:
7594: class SPCQualityGates:
7595:     """Quality validation gates for Smart Policy Chunks ingestion."""
7596:
7597:     # Phase-one output quality thresholds
7598:     MIN_CHUNKS = 5
7599:     MAX_CHUNKS = 200
7600:     MIN_CHUNK_LENGTH = 50 # characters
7601:     MAX_CHUNK_LENGTH = 5000
7602:     MIN_STRATEGIC_SCORE = 0.3
7603:     MIN_QUALITY_SCORE = 0.5
7604:     REQUIRED_CHUNK_FIELDS = ['text', 'chunk_id', 'strategic_importance', 'quality_score']
7605:
7606:     # Compatibility thresholds for downstream phases
7607:     MIN_EMBEDDING_DIM = 384 # For semantic analysis
7608:     REQUIRED_METADATA_FIELDS = ['document_id', 'title', 'version']
7609:
7610:     # CRITICAL quality metrics (per README specifications)
7611:     MIN_PROVENANCE_COMPLETENESS = 1.0 # 100% REQUIRED (no partial coverage tolerated)
7612:     MIN_STRUCTURAL_CONSISTENCY = 1.0 # 100% REQUIRED (perfect structure)
7613:     MIN_BOUNDARY_F1 = 0.85 # Chunk boundary quality
7614:     MIN_BUDGET_CONSISTENCY = 0.95 # Budget data consistency
7615:     MIN_TEMPORAL_ROBUSTNESS = 0.80 # Temporal data quality
7616:
7617:     @calibrated_method("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_input")
7618:     def validate_input(self, document_path: Path) -> dict[str, Any]:
7619:         """
7620:             Validate input document before processing.
7621:
7622:             Args:
7623:                 document_path: Path to input document
7624:
7625:             Returns:
7626:                 Dictionary with validation results
7627:         """
7628:         failures = []
7629:
7630:         # Check file exists
7631:         if not document_path.exists():
7632:             failures.append(f"Input document not found: {document_path}")
7633:             return {"passed": False, "failures": failures}
7634:
7635:         # Check file size (not empty, not too large)
7636:         file_size = document_path.stat().st_size
7637:         if file_size == 0:
7638:             failures.append("Input document is empty")
```

```
7639:     elif file_size > 50 * 1024 * 1024: # 50MB limit
7640:         failures.append(f"Input document too large: {file_size / 1024 / 1024:.1f}MB")
7641:
7642:     # Check file extension
7643:     if document_path.suffix.lower() not in ['.txt', '.pdf', '.json']:
7644:         failures.append(f"Unsupported file type: {document_path.suffix}")
7645:
7646:     return {
7647:         "passed": len(failures) == 0,
7648:         "failures": failures,
7649:         "file_size_bytes": file_size,
7650:     }
7651:
7652: @calibrated_method("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_chunks")
7653: def validate_chunks(self, chunks: list[dict[str, Any]]) -> dict[str, Any]:
7654:     """
7655:     Validate processed chunks from phase-one.
7656:
7657:     Args:
7658:         chunks: List of smart policy chunks
7659:
7660:     Returns:
7661:         Dictionary with validation results
7662:     """
7663:     failures = []
7664:     warnings = []
7665:
7666:     # Check chunk count
7667:     if len(chunks) < self.MIN_CHUNKS:
7668:         failures.append(f"Too few chunks: {len(chunks)} < {self.MIN_CHUNKS}")
7669:     elif len(chunks) > self.MAX_CHUNKS:
7670:         warnings.append(f"High chunk count: {len(chunks)} > {self.MAX_CHUNKS}")
7671:
7672:     # Validate each chunk
7673:     for idx, chunk in enumerate(chunks):
7674:         chunk_id = chunk.get('chunk_id', f'chunk_{idx}')
7675:
7676:         # Check required fields
7677:         for field in self.REQUIRED_CHUNK_FIELDS:
7678:             if field not in chunk:
7679:                 failures.append(f"{chunk_id}: Missing required field '{field}'")
7680:
7681:         # Check chunk text length
7682:         text = chunk.get('text', '')
7683:         if len(text) < self.MIN_CHUNK_LENGTH:
7684:             failures.append(f"{chunk_id}: Text too short: {len(text)} < {self.MIN_CHUNK_LENGTH}")
7685:         elif len(text) > self.MAX_CHUNK_LENGTH:
7686:             warnings.append(f"{chunk_id}: Text very long: {len(text)} > {self.MAX_CHUNK_LENGTH}")
7687:
7688:         # Check quality scores
7689:         strategic_score = chunk.get('strategic_importance', 0)
7690:         if strategic_score < self.MIN_STRATEGIC_SCORE:
7691:             warnings.append(f"{chunk_id}: Low strategic importance: {strategic_score}")
7692:
7693:         quality_score = chunk.get('quality_score', 0)
7694:         if quality_score < self.MIN_QUALITY_SCORE:
```

```
7695:             warnings.append(f"chunk_id): Low quality score: {quality_score}")
7696:
7697:         return {
7698:             "passed": len(failures) == 0,
7699:             "failures": failures,
7700:             "warnings": warnings,
7701:             "chunk_count": len(chunks),
7702:         }
7703:
7704:     @calibrated_method("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_output_compatibility")
7705:     def validate_output_compatibility(self, output: dict[str, Any]) -> dict[str, Any]:
7706:         """
7707:             Validate output structure for compatibility with downstream phases.
7708:
7709:             Ensures the phase-one output can be consumed by next phases in the flux.
7710:
7711:             Args:
7712:                 output: Phase-one output dictionary
7713:
7714:             Returns:
7715:                 Dictionary with validation results
7716:             """
7717:         failures = []
7718:
7719:         # Check required top-level keys
7720:         if 'chunks' not in output:
7721:             failures.append("Missing 'chunks' in output")
7722:
7723:         if 'metadata' not in output:
7724:             failures.append("Missing 'metadata' in output")
7725:         else:
7726:             # Validate metadata fields
7727:             metadata = output['metadata']
7728:             for field in self.REQUIRED_METADATA_FIELDS:
7729:                 if field not in metadata:
7730:                     failures.append(f"Missing required metadata field: '{field}'")
7731:
7732:         # Check chunks structure
7733:         if 'chunks' in output:
7734:             chunks_result = self.validate_chunks(output['chunks'])
7735:             if not chunks_result['passed']:
7736:                 failures.extend(chunks_result['failures'])
7737:
7738:         return {
7739:             "passed": len(failures) == 0,
7740:             "failures": failures,
7741:         }
7742:
7743:     @calibrated_method("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics")
7744:     def validate_quality_metrics(self, quality_metrics: Any) -> dict[str, Any]:
7745:         """
7746:             Validate quality metrics from CanonPolicyPackage against MAXIMUM STANDARDS.
7747:
7748:             Enforces strict thresholds per README specifications. No degradation tolerated.
7749:
7750:             Args:
```

```

7751:             quality_metrics: QualityMetrics instance from CanonPolicyPackage
7752:
7753:     Returns:
7754:         Dictionary with validation results:
7755:         {
7756:             "passed": bool,
7757:             "failures": List[str], # CRITICAL failures that MUST be fixed
7758:             "warnings": List[str], # Non-critical warnings
7759:             "metrics": Dict[str, float] # Actual metric values
7760:         }
7761:         """
7762:         failures = []
7763:         warnings = []
7764:         metrics_dict = {}
7765:
7766:         # Extract metrics (handle both object and dict)
7767:         if hasattr(quality_metrics, '__dict__'):
7768:             # It's an object
7769:                 provenance_completeness = getattr(quality_metrics, 'provenance_completeness', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L193_90", 0.0))
7770:                 structural_consistency = getattr(quality_metrics, 'structural_consistency', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L194_88", 0.0))
7771:                 boundary_f1 = getattr(quality_metrics, 'boundary_f1', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L195_66", 0.0))
7772:                 budget_consistency = getattr(quality_metrics, 'budget_consistency_score', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L196_86", 0.0))
7773:                 temporal_robustness = getattr(quality_metrics, 'temporal_robustness', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L197_82", 0.0))
7774:                 chunk_context_coverage = getattr(quality_metrics, 'chunk_context_coverage', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L198_88", 0.0))
7775:             else:
7776:                 # It's a dict
7777:                     provenance_completeness = quality_metrics.get('provenance_completeness', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L201_85", 0.0))
7778:                     structural_consistency = quality_metrics.get('structural_consistency', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L202_83", 0.0))
7779:                     boundary_f1 = quality_metrics.get('boundary_f1', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L203_61", 0.0))
7780:                     budget_consistency = quality_metrics.get('budget_consistency_score', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L204_81", 0.0))
7781:                     temporal_robustness = quality_metrics.get('temporal_robustness', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L205_77", 0.0))
7782:                     chunk_context_coverage = quality_metrics.get('chunk_context_coverage', ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "auto_param_L206_83", 0.0))
7783:
7784:             # Store actual values
7785:             metrics_dict = {
7786:                 'provenance_completeness': provenance_completeness,
7787:                 'structural_consistency': structural_consistency,
7788:                 'boundary_f1': boundary_f1,
7789:                 'budget_consistency_score': budget_consistency,
7790:                 'temporal_robustness': temporal_robustness,
7791:                 'chunk_context_coverage': chunk_context_coverage,
7792:             }
7793:
7794:             # CRITICAL: Provenance completeness MUST be 100%

```

```
7795:         if provenance_completeness < self.MIN_PROVENANCE_COMPLETENESS:
7796:             failures.append(
7797:                 f"\u03d7\237\224' CRITICAL: Provenance completeness below threshold: "
7798:                 f"\u03d7\{provenance_completeness:.2%} < {self.MIN_PROVENANCE_COMPLETENESS:.0%}. "
7799:                 f"\"Every token must be traceable to source (README requirement).\""
7800:             )
7801:             logger.error(
7802:                 f"Provenance completeness FAILED: {provenance_completeness:.2%} "
7803:                 f"(required: {self.MIN_PROVENANCE_COMPLETENESS:.0%})"
7804:             )
7805:
7806: # CRITICAL: Structural consistency MUST be perfect
7807: if structural_consistency < self.MIN_STRUCTURAL_CONSISTENCY:
7808:     failures.append(
7809:         f"\u03d7\237\224' CRITICAL: Structural consistency below threshold: "
7810:         f"\u03d7\{structural_consistency:.2%} < {self.MIN_STRUCTURAL_CONSISTENCY:.0%}. "
7811:         f"\"Policy structure must be perfectly parsed (FASE 3 gate).\""
7812:     )
7813:     logger.error(
7814:         f"Structural consistency FAILED: {structural_consistency:.2%} "
7815:         f"(required: {self.MIN_STRUCTURAL_CONSISTENCY:.0%})"
7816:     )
7817:
7818: # HIGH: Boundary F1 for chunk quality
7819: if boundary_f1 < self.MIN_BOUNDARY_F1:
7820:     failures.append(
7821:         f"\u03d7\237\224' HIGH: Boundary F1 below threshold: "
7822:         f"\u03d7\{boundary_f1:.2f} < {self.MIN_BOUNDARY_F1}. "
7823:         f"\"Chunk boundaries are not accurate enough (FASE 8 gate).\""
7824:     )
7825:     logger.error(
7826:         f"Boundary F1 FAILED: {boundary_f1:.2f} "
7827:         f"(required: {self.MIN_BOUNDARY_F1})"
7828:     )
7829:
7830: # HIGH: Budget consistency for financial data
7831: if budget_consistency < self.MIN_BUDGET_CONSISTENCY:
7832:     warnings.append(
7833:         f"\u03d7\237\237; Budget consistency below threshold: "
7834:         f"\u03d7\{budget_consistency:.2%} < {self.MIN_BUDGET_CONSISTENCY:.0%}. "
7835:         f"\"Budget data may have inconsistencies (FASE 6 gate).\""
7836:     )
7837:     logger.warning(
7838:         f"Budget consistency WARNING: {budget_consistency:.2%} "
7839:         f"(recommended: {self.MIN_BUDGET_CONSISTENCY:.0%})"
7840:     )
7841:
7842: # MEDIUM: Temporal robustness
7843: if temporal_robustness < self.MIN_TEMPORAL_ROBUSTNESS:
7844:     warnings.append(
7845:         f"\u03d7\237\237; Temporal robustness below threshold: "
7846:         f"\u03d7\{temporal_robustness:.2%} < {self.MIN_TEMPORAL_ROBUSTNESS:.0%}. "
7847:         f"\"Temporal data may be incomplete.\""
7848:     )
7849:
7850: # INFO: Chunk context coverage
```

```
7851:         if chunk_context_coverage < ParameterLoaderV2.get("farfan_core.processing.spc_ingestion.quality_gates.SPCQualityGates.validate_quality_metrics", "au
to_param_L275_36", 0.5):
7852:             warnings.append(
7853:                 f"\u2041\u217 Low chunk context coverage: {chunk_context_coverage:.2%}. "
7854:                 f"Few inter-chunk relationships detected."
7855:             )
7856:
7857:     # Summary logging
7858:     if failures:
7859:         logger.error(f"Quality metrics validation FAILED with {len(failures)} critical issues")
7860:     elif warnings:
7861:         logger.warning(f"Quality metrics validation PASSED with {len(warnings)} warnings")
7862:     else:
7863:         logger.info("Quality metrics validation PASSED - All thresholds met")
7864:
7865:     return {
7866:         "passed": len(failures) == 0,
7867:         "failures": failures,
7868:         "warnings": warnings,
7869:         "metrics": metrics_dict,
7870:     }
7871:
7872:
7873: # Legacy alias for backwards compatibility
7874: class QualityGates(SPCQualityGates):
7875:     """Legacy alias for SPCQualityGates."""
7876:     pass
7877:
7878:
7879:
7880: =====
7881: FILE: src/farfan_pipeline/processing/semantic_chunking_policy.py
7882: =====
7883:
7884: """
7885: INTERNAL SPC COMPONENT
7886:
7887: \u232 i,\u217 USAGE RESTRICTION \u232 i,\u217
7888: =====
7889: This module implements SOTA semantic chunking and policy analysis for Smart
7890: Policy Chunks. It MUST NOT be used as a standalone ingestion pipeline in the
7891: canonical FARFAN flow.
7892:
7893: Canonical entrypoint is scripts/run_policy_pipeline_verified.py.
7894:
7895: This module is an INTERNAL COMPONENT of:
7896:     src/farfan_core/processing/spc_ingestion.py (StrategicChunkingSystem)
7897:
7898: DO NOT use this module directly as an independent pipeline. It is consumed
7899: internally by the SPC core and should only be imported from within:
7900:     - farfan_core/processing/spc_ingestion
7901:     - Unit tests for SPC components
7902:
7903: Scientific Foundation:
7904: - Semantic: BGE-M3 (2024, SOTA multilingual dense retrieval)
7905: - Chunking: Semantic-aware with policy structure recognition
```

```
7906: - Math: Information-theoretic Bayesian evidence accumulation
7907: - Causal: Directed Acyclic Graph inference with interventional calculus
7908: =====
7909: """
7910: from __future__ import annotations
7911:
7912: import json
7913: import logging
7914: import re
7915: from dataclasses import dataclass
7916: from enum import Enum
7917: from typing import TYPE_CHECKING, Any, Literal
7918:
7919: import numpy as np
7920: import torch
7921: from scipy import stats
7922: from scipy.spatial.distance import cosine
7923: from scipy.special import rel_entr
7924:
7925: # Check dependency lockdown before importing transformers
7926: from farfan_pipeline.core.dependency_lockdown import get_dependency_lockdown
7927: from transformers import AutoModel, AutoTokenizer
7928: from farfan_pipeline.core.parameters import ParameterLoaderV2
7929: from farfan_pipeline.core.calibration.decorators import calibrated_method
7930:
7931: _lockdown = get_dependency_lockdown()
7932:
7933: if TYPE_CHECKING:
7934:     from numpy.typing import NDArray
7935:
7936: # Note: logging.basicConfig should be called by the application entry point,
7937: # not at module import time to avoid side effects
7938: logger = logging.getLogger("policy_framework")
7939:
7940: def _get_chunk_content(chunk: dict[str, Any]) -> str:
7941:     """Compatibility helper returning the canonical chunk content field."""
7942:
7943:     if "content" in chunk:
7944:         return chunk["content"]
7945:     return chunk.get("text", "")
7946:
7947: def _upgrade_chunk_schema(chunk: dict[str, Any]) -> dict[str, Any]:
7948:     """Return a chunk dict that guarantees ``content`` availability."""
7949:
7950:     if "content" in chunk:
7951:         return chunk
7952:     upgraded = dict(chunk)
7953:     upgraded["content"] = upgraded.get("text", "")
7954:     return upgraded
7955:
7956: # =====
7957: # CALIBRATED CONSTANTS (SOTA)
7958: # =====
7959: POSITION_WEIGHT_SCALE: float = 0.42 # Early sections exert stronger evidentiary leverage
7960: TABLE_WEIGHT_FACTOR: float = 1.35 # Tabular content is typically audited data
7961: NUMERICAL_WEIGHT_FACTOR: float = 1.18 # Numerical narratives reinforce credibility
```

```

7962: PLAN_SECTION_WEIGHT_FACTOR: float = 1.25 # Investment plans anchor execution feasibility
7963: DIAGNOSTIC_SECTION_WEIGHT_FACTOR: float = 0.92 # Diagnostics contextualize but do not commit resources
7964: RENYI_ALPHA_ORDER: float = 1.45 # Van Erven & Harremoës (2014) Optimum between KL and Rényi regimes
7965: RENYI_ALERT_THRESHOLD: float = 0.24 # Empirically tuned on 2021-2024 Colombian PDM corpus
7966: RENYI_CURVATURE_GAIN: float = 0.85 # Amplifies curvature impact without destabilizing evidence
7967: RENYI_FLUX_TEMPERATURE: float = 0.65 # Controls saturation of Renyi coherence flux
7968: RENYI_STABILITY_EPSILON: float = 1e-9 # Numerical guard-rail for degenerative posteriors
7969:
7970: # =====
7971: # DOMAIN ONTOLOGY
7972: # =====
7973:
7974: class CausalDimension(Enum):
7975:     """Marco Lógico standard (DNP Colombia)"""
7976:     INSUMOS = "insumos" # Recursos, capacidad institucional
7977:     ACTIVIDADES = "actividades" # Acciones, procesos, cronogramas
7978:     PRODUCTOS = "productos" # Entregables inmediatos
7979:     RESULTADOS = "resultados" # Efectos mediano plazo
7980:     IMPACTOS = "impactos" # Transformación estructural largo plazo
7981:     SUPUESTOS = "supuestos" # Condiciones habilitantes
7982:
7983: class PDMSection(Enum):
7984:     """
7985:         Enumerates the typical sections of a Colombian Municipal Development Plan (PDM),
7986:         as defined by Ley 152/1994. Each member represents a key structural component
7987:         of the PDM document, facilitating semantic analysis and policy structure recognition.
7988:     """
7989:     DIAGNOSTICO = "diagnóstico"
7990:     VISION_ESTRATEGICA = "visión_estratégica"
7991:     PLAN_PLURIANUAL = "plan_pluriannual"
7992:     PLAN_INVERSIONES = "plan_inversiones"
7993:     MARCO_FISCAL = "marco_fiscal"
7994:     SEGUIMIENTO = "seguimiento_evaluación"
7995:
7996: @dataclass(frozen=True, slots=True)
7997: class SemanticConfig:
7998:     """Configuración calibrada para análisis de políticas públicas"""
7999:     # BGE-M3: Best multilingual embedding (Jan 2024, beats E5)
8000:     embedding_model: str = "BAAI/bge-m3"
8001:     chunk_size: int = 768 # Optimal for policy paragraphs (empirical)
8002:     chunk_overlap: int = 128 # Preserve cross-boundary context
8003:     similarity_threshold: float = 0.82 # Calibrated on PDM corpus
8004:     min_evidence_chunks: int = 3 # Statistical significance floor
8005:     bayesian_prior_strength: float = 0.5 # Conservative uncertainty
8006:     device: Literal["cpu", "cuda"] | None = None
8007:     batch_size: int = 32
8008:     fp16: bool = True # Memory optimization
8009:
8010: # =====
8011: # SEMANTIC PROCESSOR (SOTA)
8012: # =====
8013:
8014: class SemanticProcessor:
8015:     """
8016:         State-of-the-art semantic processing with:
8017:             - BGE-M3 embeddings (2024 SOTA)

```

```

8018:     - Policy-aware chunking (respects PDM structure)
8019:     - Efficient batching with FP16
8020:     """
8021:
8022:     def __init__(self, config: SemanticConfig) -> None:
8023:         self.config = config
8024:         self._model = None
8025:         self._tokenizer = None
8026:         self._loaded = False
8027:
8028:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor._lazy_load")
8029:     def _lazy_load(self) -> None:
8030:         if self._loaded:
8031:             return
8032:         try:
8033:             device = self.config.device or ("cuda" if torch.cuda.is_available() else "cpu")
8034:             logger.info(f"Loading BGE-M3 model on {device}...")
8035:             self._tokenizer = AutoTokenizer.from_pretrained(self.config.embedding_model)
8036:             self._model = AutoModel.from_pretrained(
8037:                 self.config.embedding_model,
8038:                 torch_dtype=torch.float16 if self.config.fp16 and device == "cuda" else torch.float32
8039:             ).to(device)
8040:             self._model.eval()
8041:             self._loaded = True
8042:             logger.info("BGE-M3 loaded successfully")
8043:         except ImportError as e:
8044:             missing = None
8045:             msg = str(e)
8046:             if "transformers" in msg:
8047:                 missing = "transformers"
8048:             elif "torch" in msg:
8049:                 missing = "torch"
8050:             else:
8051:                 missing = "transformers or torch"
8052:             raise RuntimeError(
8053:                 f"Missing dependency: {missing}. Please install with 'pip install {missing}'"
8054:             ) from e
8055:
8056:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor.chunk_text")
8057:     def chunk_text(self, text: str, preserve_structure: bool = True) -> list[dict[str, Any]]:
8058:         """
8059:             Policy-aware semantic chunking:
8060:             - Respects section boundaries (numbered lists, headers)
8061:             - Maintains table integrity
8062:             - Preserves reference links between text segments
8063:         """
8064:         self._lazy_load()
8065:         # Detect structural elements (headings, numbered sections, tables)
8066:         if preserve_structure:
8067:             sections = self._detect_pdm_structure(text)
8068:         else:
8069:             sections = [{"text": text, "type": "TEXT", "id": 0}]
8070:         chunks = []
8071:         for section in sections:
8072:             # Tokenize section
8073:             tokens = self._tokenizer.encode(

```

```

8074:             section["text"],
8075:             add_special_tokens=False,
8076:             truncation=False
8077:         )
8078:         # Sliding window with overlap
8079:         for i in range(0, len(tokens), self.config.chunk_size - self.config.chunk_overlap):
8080:             chunk_tokens = tokens[i:i + self.config.chunk_size]
8081:             chunk_text = self._tokenizer.decode(chunk_tokens, skip_special_tokens=True)
8082:             chunks.append({
8083:                 "content": chunk_text,
8084:                 "section_type": section["type"],
8085:                 "section_id": section["id"],
8086:                 "token_count": len(chunk_tokens),
8087:                 "position": len(chunks),
8088:                 "has_table": self._detect_table(chunk_text),
8089:                 "has_numerical": self._detect_numerical_data(chunk_text),
8090:                 "pdq_context": {},
8091:             })
8092:         # Batch embed all chunks
8093:         embeddings = self._embed_batch([c["content"] for c in chunks])
8094:         for chunk, emb in zip(chunks, embeddings, strict=False):
8095:             chunk["embedding"] = emb
8096:         logger.info(f"Generated {len(chunks)} policy-aware chunks")
8097:         return [_upgrade_chunk_schema(chunk) for chunk in chunks]
8098:
8099:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor._detect_pdm_structure")
8100:     def _detect_pdm_structure(self, text: str) -> list[dict[str, Any]]:
8101:         """Detect PDM sections using Colombian policy document patterns"""
8102:         sections = []
8103:         # Patterns for Colombian PDM structure
8104:         patterns = {
8105:             PDMSection.DIAGNOSTICO: r"(?i)(diagnÃ³stico|caracterizaciÃ³n|situaciÃ³n actual)",
8106:             PDMSection.VISION_ESTRATEGICA: r"(?i)(visiÃ³n|misiÃ³n|objetivos estratÃ©gicos)",
8107:             PDMSection.PLURIANUAL: r"(?i)(plan plurianual|programas|proyectos)",
8108:             PDMSection.PLAN_INVERSIONES: r"(?i)(plan de inversiones|presupuesto|recursos)",
8109:             PDMSection.MARCO_FISCAL: r"(?i)(marco fiscal|sostenibilidad fiscal)",
8110:             PDMSection.SEGUIMIENTO: r"(?i)(seguimiento|evaluaciÃ³n|indicadores)"
8111:         }
8112:         # Split by major headers (numbered or capitalized)
8113:         parts = re.split(r'\n(?=[0-9]+\.|[A-ZÀ\221À\201À\211À\215À\223À\232]{3,})', text)
8114:         for i, part in enumerate(parts):
8115:             section_type = PDMSection.DIAGNOSTICO # default
8116:             for stype, pattern in patterns.items():
8117:                 if re.search(pattern, part[:200]):
8118:                     section_type = stype
8119:                     break
8120:             sections.append({
8121:                 "text": part.strip(),
8122:                 "type": section_type,
8123:                 "id": f"sec_{i}"
8124:             })
8125:         return sections
8126:
8127:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor._detect_table")
8128:     def _detect_table(self, text: str) -> bool:
8129:         """Detect if chunk contains tabular data"""

```

```

8130:     # Multiple tabs or pipes suggest table structure
8131:     return (text.count('\t') > 3 or
8132:             text.count(' | ') > 3 or
8133:             bool(re.search(r'\d+\s+\d+\s+\d+', text)))
8134:
8135: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor._detect_numerical_data")
8136: def _detect_numerical_data(self, text: str) -> bool:
8137:     """Detect if chunk contains significant numerical/financial data"""
8138:     # Look for currency, percentages, large numbers
8139:     patterns = [
8140:         r'${}\d+(?:[.,]\d+)*',    # Currency
8141:         r'\d+(?:[.,]\d+)*\s*%',   # Percentages
8142:         r'\d{1,3}(?:[.,]\d{3})+', # Large numbers with separators
8143:     ]
8144:     return any(re.search(p, text) for p in patterns)
8145:
8146: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor._embed_batch")
8147: def _embed_batch(self, texts: list[str]) -> list[NDArray[np.floating[Any]]]:
8148:     """Batch embedding with BGE-M3"""
8149:     self._lazy_load()
8150:     embeddings = []
8151:     for i in range(0, len(texts), self.config.batch_size):
8152:         batch = texts[i:i + self.config.batch_size]
8153:         # Tokenize batch
8154:         encoded = self._tokenizer(
8155:             batch,
8156:             padding=True,
8157:             truncation=True,
8158:             max_length=self.config.chunk_size,
8159:             return_tensors="pt"
8160:         ).to(self._model.device)
8161:         # Generate embeddings (mean pooling)
8162:         with torch.no_grad():
8163:             outputs = self._model(**encoded)
8164:             # Mean pooling over sequence
8165:             attention_mask = encoded["attention_mask"]
8166:             token_embeddings = outputs.last_hidden_state
8167:             input_mask_expanded = attention_mask.unsqueeze(-1).expand(token_embeddings.size()).float()
8168:             sum_embeddings = torch.sum(token_embeddings * input_mask_expanded, 1)
8169:             sum_mask = torch.clamp(input_mask_expanded.sum(1), min=1e-9)
8170:             batch_embeddings = (sum_embeddings / sum_mask).cpu().numpy()
8171:             embeddings.extend([emb.astype(np.float32) for emb in batch_embeddings])
8172:     return embeddings
8173:
8174: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticProcessor.embed_single")
8175: def embed_single(self, text: str) -> NDArray[np.floating[Any]]:
8176:     """Single text embedding"""
8177:     return self._embed_batch([text])[0]
8178:
8179: # =====
8180: # MATHEMATICAL ENHANCER (RIGOROUS)
8181: # =====
8182:
8183: class BayesianEvidenceIntegrator:
8184:     """
8185:         Information-theoretic Bayesian evidence accumulation:

```

```

8186:     - Dirichlet-Multinomial for multi-hypothesis tracking
8187:     - KL divergence for belief update quantification
8188:     - Entropy-based confidence calibration
8189:     - No simplifications or heuristics
8190:     """
8191:
8192:     def __init__(self, prior_concentration: float = 0.5) -> None:
8193:         """
8194:             Args:
8195:                 prior_concentration: Dirichlet concentration ( $\hat{t} \pm$ ).
8196:                     Lower = more uncertain prior (conservative)
8197:             """
8198:             if prior_concentration <= 0:
8199:                 raise ValueError(
8200:                     f"Invalid prior_concentration: Dirichlet concentration parameter must be strictly positive. "
8201:                     f"Received: {prior_concentration}"
8202:                 )
8203:             self.prior_alpha = float(prior_concentration)
8204:
8205:             def integrate_evidence(
8206:                 self,
8207:                 similarities: NDArray[np.float64],
8208:                 chunk_metadata: list[dict[str, Any]]
8209:             ) -> dict[str, float]:
8210:                 """
8211:                     Bayesian evidence integration with information-theoretic rigor:
8212:                     1. Map similarities to likelihood space via monotonic transform
8213:                     2. Weight evidence by chunk reliability (position, structure, content type)
8214:                     3. Update Dirichlet posterior
8215:                     4. Compute information gain (KL divergence from prior)
8216:                     5. Calculate calibrated confidence with epistemic uncertainty
8217:                 """
8218:                 if len(similarities) == 0:
8219:                     return self._null_evidence()
8220:                 # 1. Transform similarities to probability space
8221:                 # Using sigmoid with learned temperature for calibration
8222:                 sims = np.asarray(similarities, dtype=np.float64)
8223:                 probs = self._similarity_to_probability(sims)
8224:                 # 2. Compute reliability weights from metadata
8225:                 weights = self._compute_reliability_weights(chunk_metadata)
8226:                 # 3. Aggregate weighted evidence
8227:                 # Dirichlet posterior parameters:  $\hat{t}_{\text{post}} = \hat{t}_{\text{prior}} + \text{weighted\_counts}$ 
8228:                 positive_evidence = np.sum(weights * probs)
8229:                 negative_evidence = np.sum(weights * (ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator.__init__", "auto_param_L348_46", 1.0) - probs))
8230:                 alpha_pos = self.prior_alpha + positive_evidence
8231:                 alpha_neg = self.prior_alpha + negative_evidence
8232:                 alpha_total = alpha_pos + alpha_neg
8233:                 # 4. Posterior statistics
8234:                 posterior_mean = alpha_pos / alpha_total
8235:                 posterior_variance = (alpha_pos * alpha_neg) / (
8236:                     alpha_total**2 * (alpha_total + 1)
8237:                 )
8238:                 # 5. Information gain (KL divergence from prior to posterior)
8239:                 prior_dist = np.array([self.prior_alpha, self.prior_alpha])
8240:                 prior_dist = prior_dist / prior_dist.sum()

```

```

8241:     posterior_dist = np.array([alpha_pos, alpha_neg])
8242:     posterior_dist = posterior_dist / posterior_dist.sum()
8243:     kl_divergence = float(np.sum(rel_entr(posterior_dist, prior_dist)))
8244:     # 6. Entropy-based calibrated confidence
8245:     posterior_entropy = stats.beta.entropy(alpha_pos, alpha_neg)
8246:     max_entropy = stats.beta.entropy(1, 1) # Maximum uncertainty
8247:     confidence = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator.__init__", "auto_param_L366_21", 1.0)
8248:     return {
8249:         "posterior_mean": float(np.clip(posterior_mean, ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator.__init__", "auto_param_L368_60", 0.0), ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator.__init__", "auto_param_L368_65", 1.0))),
8250:         "posterior_std": float(np.sqrt(posterior_variance)),
8251:         "information_gain": float(kl_divergence),
8252:         "confidence": float(confidence),
8253:         "evidence_strength": float(
8254:             positive_evidence / (alpha_total - 2 * self.prior_alpha)
8255:             if abs(alpha_total - 2 * self.prior_alpha) > 1e-8 else ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator.__init__", "auto_param_L374_71", 0.0)
8256:         ),
8257:         "n_chunks": len(similarities)
8258:     }
8259:
8260:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._similarity_to_probability")
8261:     def _similarity_to_probability(self, sims: NDArray[np.float64]) -> NDArray[np.float64]:
8262:         """
8263:             Calibrated transform from cosine similarity [-1,1] to probability [0,1]
8264:             Using sigmoid with empirically derived temperature
8265:         """
8266:         # Shift to [0,2], scale to reasonable range
8267:         x = (sims + ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._similarity_to_probability", "auto_param_L386_20", 1.0)) * 2.0
8268:         # Sigmoid with temperature=2.0 (calibrated on policy corpus)
8269:         return ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._similarity_to_probability", "auto_param_L386_15", 1.0) / (ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._similarity_to_probability", "auto_param_L388_22", 1.0) + np.exp(-x / 2.0))
8270:
8271:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._compute_reliability_weights")
8272:     def _compute_reliability_weights(self, metadata: list[dict[str, Any]]) -> NDArray[np.float64]:
8273:         """
8274:             Evidence reliability based on:
8275:             - Position in document (early sections more diagnostic)
8276:             - Content type (tables/numbers more reliable for quantitative claims)
8277:             - Section type (plan sections more reliable than diagnostics)
8278:         """
8279:         n = len(metadata)
8280:         weights = np.ones(n, dtype=np.float64)
8281:         for i, meta in enumerate(metadata):
8282:             # Position weight (early = more reliable)
8283:             pos_weight = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._compute_reliability_weights", "auto_param_L402_25", 1.0) - (meta["position"] / max(1, n)) * POSITION_WEIGHT_SCALE
8284:             # Content type weight
8285:             content_weight = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._compute_reliability_weights", "content_weight", 1.0) # Refactored
8286:             if meta.get("has_table", False):
8287:                 content_weight *= TABLE_WEIGHT_FACTOR

```

```

8288:         if meta.get("has_numerical", False):
8289:             content_weight *= NUMERICAL_WEIGHT_FACTOR
8290:             # Section type weight (plan sections > diagnostic)
8291:             section_type = meta.get("section_type")
8292:             if section_type in [PDMSection.PLAN_PLURIANUAL, PDMSection.PLAN_INVERSIONES]:
8293:                 content_weight *= PLAN_SECTION_WEIGHT_FACTOR
8294:             elif section_type == PDMSection.DIAGNOSTICO:
8295:                 content_weight *= DIAGNOSTIC_SECTION_WEIGHT_FACTOR
8296:             weights[i] = pos_weight * content_weight
8297:             # Normalize to sum to n (preserve total evidence mass)
8298:             return weights * (n / weights.sum())
8299:
8300:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence")
8301:     def _null_evidence(self) -> dict[str, float]:
8302:         """Return prior state (no evidence)"""
8303:         prior_mean = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "prior_mean", 0.5) #
Refactored
8304:         prior_var = self.prior_alpha / \
8305:             ((2 * self.prior_alpha)**2 * (2 * self.prior_alpha + 1))
8306:         return {
8307:             "posterior_mean": prior_mean,
8308:             "posterior_std": float(np.sqrt(prior_var)),
8309:             "information_gain": ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L428_32", 0.0),
8310:             "confidence": ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L426_0.0"),
8311:             "evidence_strength": ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L430_33", 0.0),
8312:             "n_chunks": 0
8313:         }
8314:
8315:     def causal_strength(
8316:         self,
8317:         cause_emb: NDArray[np.floating[Any]],
8318:         effect_emb: NDArray[np.floating[Any]],
8319:         context_emb: NDArray[np.floating[Any]]
8320:     ) -> float:
8321:         """
8322:             Causal strength via conditional independence approximation:
8323:             strength = sim(cause, effect) * [1 - |sim(cause, ctx) - sim(effect, ctx)|]
8324:             Intuition: Strong causal link if cause-effect similar AND
8325:             both relate similarly to context (conditional independence test proxy)
8326:             """
8327:             sim_ce = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L446_17", 1.0) - cosine(cause_emb, effect_emb)
8328:             sim_c_ctx = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L447_20", 1.0) - cosine(cause_emb, context_emb)
8329:             sim_e_ctx = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L448_20", 1.0) - cosine(effect_emb, context_emb)
8330:             # Conditional independence proxy
8331:             cond_indep = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L450_21", 1.0) - abs(sim_c_ctx - sim_e_ctx)
8332:             # Combined strength (normalized to [0,1])
8333:             strength = ((sim_ce + 1) / 2) * cond_indep
8334:             return float(np.clip(strength, ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L453_39", 0.0), ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.BayesianEvidenceIntegrator._null_evidence", "auto_param_L453_44", 1.0)))

```

```
1.0)))
8335:
8336: # =====
8337: # POLICY ANALYZER (INTEGRATED)
8338: # =====
8339:
8340: class PolicyDocumentAnalyzer:
8341:     """
8342:         Colombian Municipal Development Plan Analyzer:
8343:             - BGE-M3 semantic processing
8344:             - Policy-aware chunking (respects PDM structure)
8345:             - Bayesian evidence integration with information theory
8346:             - Causal dimension analysis per Marco LÃ³gico
8347:     """
8348:
8349:     def __init__(self, config: SemanticConfig | None = None) -> None:
8350:         self.config = config or SemanticConfig()
8351:         self.semantic = SemanticProcessor(self.config)
8352:         self.bayesian = BayesianEvidenceIntegrator(
8353:             prior_concentration=self.config.bayesian_prior_strength
8354:         )
8355:         # Initialize dimension embeddings
8356:         self.dimension_embeddings = self._init_dimension_embeddings()
8357:
8358: @calibrated_method("farfan_core.processing.semantic_chunking_policy.PolicyDocumentAnalyzer._init_dimension_embeddings")
8359: def _init_dimension_embeddings(self) -> dict[CausalDimension, NDArray[np.floating[Any]]]:
8360:     """
8361:         Canonical embeddings for Marco LÃ³gico dimensions
8362:         Using Colombian policy-specific terminology
8363:     """
8364:     descriptions = {
8365:         CausalDimension.INSUMOS: (
8366:             "recursos humanos financieros tÃ©cnicos capacidad institucional "
8367:             "presupuesto asignado infraestructura disponible personal capacitado"
8368:         ),
8369:         CausalDimension.ACTIVIDADES: (
8370:             "actividades programadas acciones ejecutadas procesos implementados "
8371:             "cronograma cumplido capacitaciones realizadas gestiones adelantadas"
8372:         ),
8373:         CausalDimension.PRODUCTOS: (
8374:             "productos entregables resultados inmediatos bienes servicios generados "
8375:             "documentos producidos obras construidas beneficiarios atendidos"
8376:         ),
8377:         CausalDimension.RESULTADOS: (
8378:             "resultados efectos mediano plazo cambios comportamiento acceso mejorado "
8379:             "capacidades fortalecidas servicios prestados metas alcanzadas"
8380:         ),
8381:         CausalDimension.IMPACTOS: (
8382:             "impactos transformaciÃ³n estructural efectos largo plazo desarrollo sostenible "
8383:             "bienestar poblacional reducciÃ³n pobreza equidad territorial"
8384:         ),
8385:         CausalDimension.SUPUESTOS: (
8386:             "supuestos condiciones habilitantes riesgos externos factores contextuales "
8387:             "viabilidad polÃtica sostenibilidad financiera apropiaciÃ³n comunitaria"
8388:         )
8389:     }
```

```

8390:         return {
8391:             dim: self.semantic.embed_single(desc)
8392:             for dim, desc in descriptions.items()
8393:         }
8394:
8395:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.PolicyDocumentAnalyzer.analyze")
8396:     def analyze(self, text: str) -> dict[str, Any]:
8397:         """
8398:             Full pipeline: chunking \206\222 embedding \206\222 dimension analysis \206\222 evidence integration
8399:         """
8400:         # 1. Policy-aware chunking
8401:         chunks = self.semantic.chunk_text(text, preserve_structure=True)
8402:         logger.info(f"Processing {len(chunks)} chunks")
8403:         # 2. Analyze each causal dimension
8404:         dimension_results = {}
8405:         for dim, dim_emb in self.dimension_embeddings.items():
8406:             similarities = np.array([
8407:                 ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.PolicyDocumentAnalyzer.analyze", "auto_param_L526_16", 1.0) - cosine(
chunk["embedding"], dim_emb)
8408:                 for chunk in chunks
8409:             ])
8410:             # Filter by threshold
8411:             relevant_mask = similarities >= self.config.similarity_threshold
8412:             relevant_sims = similarities[relevant_mask]
8413:             relevant_chunks = [c for c, m in zip(chunks, relevant_mask, strict=False) if m]
8414:             # Bayesian integration
8415:             if len(relevant_sims) >= self.config.min_evidence_chunks:
8416:                 evidence = self.bayesian.integrate_evidence(
8417:                     relevant_sims,
8418:                     relevant_chunks
8419:                 )
8420:             else:
8421:                 evidence = self.bayesian._null_evidence()
8422:             dimension_results[dim.value] = {
8423:                 "total_chunks": int(np.sum(relevant_mask)),
8424:                 "mean_similarity": float(np.mean(similarities)),
8425:                 "max_similarity": float(np.max(similarities)),
8426:                 **evidence
8427:             }
8428:             # 3. Extract key findings (top chunks per dimension)
8429:             key_excerpts = self._extract_key_excerpts(chunks, dimension_results)
8430:             return {
8431:                 "summary": {
8432:                     "total_chunks": len(chunks),
8433:                     "sections_detected": len({c["section_type"] for c in chunks}),
8434:                     "has_tables": sum(1 for c in chunks if c["has_table"]),
8435:                     "has_numerical": sum(1 for c in chunks if c["has_numerical"])
8436:                 },
8437:                 "causal_dimensions": dimension_results,
8438:                 "key_excerpts": key_excerpts
8439:             }
8440:
8441:     def _extract_key_excerpts(
8442:         self,
8443:         chunks: list[dict[str, Any]],
8444:         dimension_results: dict[str, dict[str, Any]]

```

```

8445:     ) -> dict[str, list[str]]:
8446:         """Extract most relevant text excerpts per dimension"""
8447:         _ = dimension_results # parameter kept for future compatibility
8448:         excerpts = {}
8449:         for dim, dim_emb in self.dimension_embeddings.items():
8450:             # Rank chunks by similarity
8451:             sims = [
8452:                 (i, ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.PolicyDocumentAnalyzer.analyze", "auto_param_L571_20", 1.0) - cos
ine(chunk["embedding"], dim_emb))
8453:                     for i, chunk in enumerate(chunks)
8454:                 ]
8455:             sims.sort(key=lambda x: x[1], reverse=True)
8456:             # Top 3 excerpts
8457:             top_chunks = [chunks[i] for i, _ in sims[:3]]
8458:             excerpts[dim.value] = [
8459:                 _get_chunk_content(c)[:300]
8460:                 + ("..." if len(_get_chunk_content(c)) > 300 else "")
8461:                 for c in top_chunks
8462:             ]
8463:         return excerpts
8464:
8465: # =====
8466: # PRODUCER CLASS - Registry Exposure
8467: # =====
8468:
8469: class SemanticChunkingProducer:
8470:     """
8471:         Producer wrapper for semantic chunking and policy analysis with registry exposure
8472:
8473:         Provides public API methods for orchestrator integration without exposing
8474:         internal implementation details or summarization logic.
8475:
8476:         Version: ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.PolicyDocumentAnalyzer.analyze", "auto_param_L595_13", 1.0).0
8477:         Producer Type: Semantic Analysis / Chunking
8478:     """
8479:
8480:     def __init__(self, config: SemanticConfig | None = None) -> None:
8481:         """Initialize producer with optional configuration"""
8482:         self.config = config or SemanticConfig()
8483:         self.semantic = SemanticProcessor(self.config)
8484:         self.bayesian = BayesianEvidenceIntegrator(
8485:             prior_concentration=self.config.bayesian_prior_strength
8486:         )
8487:         self.analyzer = PolicyDocumentAnalyzer(self.config)
8488:         logger.info("SemanticChunkingProducer initialized")
8489:
8490: # =====
8491: # CHUNKING API
8492: # =====
8493:
8494: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.chunk_document")
8495: def chunk_document(self, text: str, preserve_structure: bool = True) -> list[dict[str, Any]]:
8496:     """Chunk document into semantic units with embeddings"""
8497:     return self.semantic.chunk_text(text, preserve_structure)
8498:
8499: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_chunk_count")

```

```
8500:     def get_chunk_count(self, chunks: list[dict[str, Any]]) -> int:
8501:         """Get number of chunks"""
8502:         return len(chunks)
8503:
8504:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_chunk_text")
8505:     def get_chunk_text(self, chunk: dict[str, Any]) -> str:
8506:         """Extract text from chunk"""
8507:         return _get_chunk_content(chunk)
8508:
8509:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_chunk_embedding")
8510:     def get_chunk_embedding(self, chunk: dict[str, Any]) -> NDArray[np.floating[Any]]:
8511:         """Extract embedding from chunk"""
8512:         return chunk.get("embedding", np.array([]))
8513:
8514:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_chunk_metadata")
8515:     def get_chunk_metadata(self, chunk: dict[str, Any]) -> dict[str, Any]:
8516:         """Extract metadata from chunk"""
8517:         return {
8518:             "section_type": chunk.get("section_type"),
8519:             "section_id": chunk.get("section_id"),
8520:             "token_count": chunk.get("token_count"),
8521:             "position": chunk.get("position"),
8522:             "has_table": chunk.get("has_table"),
8523:             "has_numerical": chunk.get("has_numerical")
8524:         }
8525:
8526:     # =====
8527:     # EMBEDDING API
8528:     # =====
8529:
8530:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.embed_text")
8531:     def embed_text(self, text: str) -> NDArray[np.floating[Any]]:
8532:         """Generate single embedding for text"""
8533:         return self.semantic.embed_single(text)
8534:
8535:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.embed_batch")
8536:     def embed_batch(self, texts: list[str]) -> list[NDArray[np.floating[Any]]]:
8537:         """Generate embeddings for batch of texts"""
8538:         return self.semantic._embed_batch(texts)
8539:
8540:     # =====
8541:     # ANALYSIS API
8542:     # =====
8543:
8544:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.analyze_document")
8545:     def analyze_document(self, text: str) -> dict[str, Any]:
8546:         """Full pipeline analysis of document"""
8547:         return self.analyzer.analyze(text)
8548:
8549:     def get_dimension_analysis(
8550:         self,
8551:         analysis: dict[str, Any],
8552:         dimension: CausalDimension
8553:     ) -> dict[str, Any]:
8554:         """Extract specific dimension results from analysis"""
8555:         return analysis.get("causal_dimensions", {}).get(dimension.value, {})
```

```
8556:
8557:     def get_dimension_score(
8558:         self,
8559:         analysis: dict[str, Any],
8560:         dimension: CausalDimension
8561:     ) -> float:
8562:         """Extract dimension evidence strength score"""
8563:         dim_result = self.get_dimension_analysis(analysis, dimension)
8564:         return dim_result.get("evidence_strength", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.analyze_document", "auto_param_L683_51", 0.0))
8565:
8566:     def get_dimension_confidence(
8567:         self,
8568:         analysis: dict[str, Any],
8569:         dimension: CausalDimension
8570:     ) -> float:
8571:         """Extract dimension confidence score"""
8572:         dim_result = self.get_dimension_analysis(analysis, dimension)
8573:         return dim_result.get("confidence", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.analyze_document", "auto_param_L692_44", 0.0))
8574:
8575:     def get_dimension_excerpts(
8576:         self,
8577:         analysis: dict[str, Any],
8578:         dimension: CausalDimension
8579:     ) -> list[str]:
8580:         """Extract key excerpts for dimension"""
8581:         return analysis.get("key_excerpt", {}).get(dimension.value, [])
8582:
8583: # =====
8584: # BAYESIAN EVIDENCE API
8585: # =====
8586:
8587:     def integrate_evidence(
8588:         self,
8589:         similarities: NDArray[np.float64],
8590:         chunk_metadata: list[dict[str, Any]]
8591:     ) -> dict[str, float]:
8592:         """Perform Bayesian evidence integration"""
8593:         return self.bayesian.integrate_evidence(similarities, chunk_metadata)
8594:
8595:     def calculate_causal_strength(
8596:         self,
8597:         cause_emb: NDArray[np.floating[Any]],
8598:         effect_emb: NDArray[np.floating[Any]],
8599:         context_emb: NDArray[np.floating[Any]]
8600:     ) -> float:
8601:         """Calculate causal strength between embeddings"""
8602:         return self.bayesian.causal_strength(cause_emb, effect_emb, context_emb)
8603:
8604: @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_posterior_mean")
8605: def get_posterior_mean(self, evidence: dict[str, float]) -> float:
8606:     """Extract posterior mean from evidence integration"""
8607:     return evidence.get("posterior_mean", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_posterior_mean", "auto_param_L726_46", 0.0))
8608:
```

```

8609:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_posterior_std")
8610:     def get_posterior_std(self, evidence: dict[str, float]) -> float:
8611:         """Extract posterior standard deviation"""
8612:         return evidence.get("posterior_std", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_posterior_s
td", "auto_param_L731_45", 0.0))
8613:
8614:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_information_gain")
8615:     def get_information_gain(self, evidence: dict[str, float]) -> float:
8616:         """Extract information gain (KL divergence)"""
8617:         return evidence.get("information_gain", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_informati
on_gain", "auto_param_L736_48", 0.0))
8618:
8619:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_confidence")
8620:     def get_confidence(self, evidence: dict[str, float]) -> float:
8621:         """Extract confidence score"""
8622:         return evidence.get("confidence", ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_confidence", "a
uto_param_L741_42", 0.0))
8623:
8624:     # =====
8625:     # SEMANTIC SEARCH API
8626:     # =====
8627:
8628:     def semantic_search(
8629:         self,
8630:         query: str,
8631:         chunks: list[dict[str, Any]],
8632:         dimension: CausalDimension | None = None,
8633:         top_k: int = 5
8634:     ) -> list[tuple[dict[str, Any], float]]:
8635:         """Search chunks semantically for query"""
8636:         query_emb = self.semantic.embed_single(query)
8637:
8638:         results = []
8639:         for chunk in chunks:
8640:             chunk_emb = chunk.get("embedding")
8641:             if chunk_emb is not None and len(chunk_emb) > 0:
8642:                 similarity = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_confidence", "auto_param_L7
61_29", 1.0) - cosine(query_emb, chunk_emb)
8643:
8644:                 # Filter by dimension if specified
8645:                 if dimension is None or chunk.get("section_type") == dimension:
8646:                     results.append((chunk, float(similarity)))
8647:
8648:         # Sort by similarity descending
8649:         results.sort(key=lambda x: x[1], reverse=True)
8650:
8651:         return results[:top_k]
8652:
8653:     # =====
8654:     # UTILITY API
8655:     # =====
8656:
8657:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.list_dimensions")
8658:     def list_dimensions(self) -> list[CausalDimension]:
8659:         """List all causal dimensions"""
8660:         return list(CausalDimension)

```

```
8661:
8662:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_dimension_description")
8663:     def get_dimension_description(self, dimension: CausalDimension) -> str:
8664:         """Get description for dimension"""
8665:         descriptions = {
8666:             CausalDimension.INSUMOS: "Recursos, capacidad institucional",
8667:             CausalDimension.ACTIVIDADES: "Acciones, procesos, cronogramas",
8668:             CausalDimension.PRODUCTOS: "Entregables inmediatos",
8669:             CausalDimension.RESULTADOS: "Efectos mediano plazo",
8670:             CausalDimension.IMPACTOS: "TransformaciÃ³n estructural largo plazo",
8671:             CausalDimension.SUPUESTOS: "Condiciones habilitantes"
8672:         }
8673:         return descriptions.get(dimension, "")
8674:
8675:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.get_config")
8676:     def get_config(self) -> SemanticConfig:
8677:         """Get current configuration"""
8678:         return self.config
8679:
8680:     @calibrated_method("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.set_config")
8681:     def set_config(self, config: SemanticConfig) -> None:
8682:         """Update configuration (requires reinitialization)"""
8683:         self.config = config
8684:         self.semantic = SemanticProcessor(self.config)
8685:         self.bayesian = BayesianEvidenceIntegrator(
8686:             prior_concentration=self.config.bayesian_prior_strength
8687:         )
8688:         self.analyzer = PolicyDocumentAnalyzer(self.config)
8689:
8690: # =====
8691: # CLI INTERFACE
8692: # =====
8693:
8694: def main() -> None:
8695:     """Example usage"""
8696:     sample_pdm = """
8697: PLAN DE DESARROLLO MUNICIPAL 2024-2027
8698: MUNICIPIO DE EJEMPLO, COLOMBIA
8699:
8700: 1. DIAGNÃ\223STICO TERRITORIAL
8701: El municipio cuenta con 45,000 habitantes, de los cuales 60% reside en zona rural.
8702: La tasa de pobreza multidimensional es 42.3%, superior al promedio departamental.
8703:
8704: 2. VISIÃ\223N ESTRATÃ\211GICA
8705: Para 2027, el municipio serÃ¡ reconocido por su desarrollo sostenible e inclusivo.
8706:
8707: 3. PLAN PLURIANUAL DE INVERSIONES
8708: Se destinarÃ¡n $12,500 millones al sector educaciÃ³n, con meta de construir
8709: 3 instituciones educativas y capacitar 250 docentes en pedagogÃ-as innovadoras.
8710:
8711: 4. SEGUIMIENTO Y EVALUACIÃ\223N
8712: Se implementarÃ¡ un sistema de indicadores alineado con ODS, con mediciones semestrales.
8713: """
8714:     config = SemanticConfig(
8715:         chunk_size=512,
8716:         chunk_overlap=100,
```

```
8717:     similarity_threshold = ParameterLoaderV2.get("farfan_core.processing.semantic_chunking_policy.SemanticChunkingProducer.set_config", "similarity_threshold", 0.8) # Refactored
8718:     )
8719:     analyzer = PolicyDocumentAnalyzer(config)
8720:     results = analyzer.analyze(sample_pdm)
8721:     print(json.dumps({
8722:         "summary": results["summary"],
8723:         "dimensions": {
8724:             k: {
8725:                 "evidence_strength": v["evidence_strength"],
8726:                 "confidence": v["confidence"],
8727:                 "information_gain": v["information_gain"]
8728:             }
8729:             for k, v in results["causal_dimensions"].items()
8730:         },
8731:     }, indent=2, ensure_ascii=False))
8732:
8733:
8734:
8735: =====
8736: FILE: src/farfan_pipeline/processing/spc_ingestion.py
8737: =====
8738:
8739: """
8740: SISTEMA INDUSTRIAL SOTA PARA SMART-POLICY-CHUNKS DE PLANES DE DESARROLLO
8741: VERSIÃN 3.0 COMPLETA - SIN PLACEHOLDERS, IMPLEMENTACIÃN TOTAL
8742: FASE 1 DEL PIPELINE: GENERACIÃN DE CHUNKS COMPRENSIVOS, RIGUROSOS Y ESTRATÃGICOS
8743: """
8744:
8745: import os
8746: import re
8747: import logging
8748: import hashlib
8749: import numpy as np
8750: import copy
8751: from dataclasses import dataclass, asdict, field
8752: from typing import Dict, List, Any, Optional, Tuple, Set, Union
8753: from enum import Enum
8754: from pathlib import Path
8755: from scipy.spatial.distance import cosine
8756: from scipy.stats import entropy
8757: from scipy.signal import find_peaks
8758: # Note: torch and transformers imports removed - model lifecycle managed by canonical producers
8759: from datetime import datetime, timezone
8760: from collections import defaultdict, Counter
8761: import json
8762: import networkx as nx
8763: from sklearn.cluster import DBSCAN, AgglomerativeClustering
8764: # Note: cosine_similarity removed - using canonical semantic_search with cross-encoder reranking
8765: from sklearn.decomposition import LatentDirichletAllocation
8766: from sklearn.feature_extraction.text import TfidfVectorizer
8767: import spacy
8768: from nltk.tokenize import sent_tokenize
8769: # Note: SentenceTransformer import removed - embedding handled by canonical producers
8770: import warnings
8771: warnings.filterwarnings('ignore')
```

```
8772:  
8773: # =====  
8774: # CANONICAL MODULE INTEGRATION - SOTA Producer APIs  
8775: # Import production-grade canonical components from farfan_core.processing  
8776: # These replace internal duplicate implementations with frontier SOTA approaches  
8777: # =====  
8778:  
8779: from farfan_pipeline.processing.embedding_policy import EmbeddingPolicyProducer  
8780: from farfan_pipeline.processing.policy_processor import create_policy_processor  
8781: from farfan_pipeline.processing.semantic_chunking_policy import SemanticChunkingProducer  
8782:  
8783: # =====  
8784: # LOGGING CONFIGURADO  
8785: # =====  
8786:  
8787: logging.basicConfig(  
8788:     level=logging.INFO,  
8789:     format='%(asctime)s - %(name)s - %(levelname)s - %(message)s',  
8790:     handlers=[  
8791:         logging.FileHandler('smart_chunks_pipeline.log'),  
8792:         logging.StreamHandler()  
8793:     ]  
8794: )  
8795: logger = logging.getLogger("SPC")  
8796:  
8797: # Optional language detection for multi-language support  
8798: try:  
8799:     from langdetect import detect  
8800:     LANGDETECT_AVAILABLE = True  
8801: except ImportError:  
8802:     LANGDETECT_AVAILABLE = False  
8803:     logger.warning("langdetect not available - defaulting to Spanish models")  
8804:  
8805: # =====  
8806: # UTILITY FUNCTIONS - Serialization, Hashing, Text Safety  
8807: # =====  
8808:  
8809: def np_to_list(obj):  
8810:     """  
8811:     Convert NumPy arrays to lists for JSON serialization.  
8812:  
8813:     Inputs:  
8814:         obj: Any Python object, typically a NumPy array  
8815:     Outputs:  
8816:         List representation of the array if input is ndarray  
8817:     Raises:  
8818:         TypeError if object is not JSON-serializable  
8819:     """  
8820:     if isinstance(obj, np.ndarray):  
8821:         return obj.tolist()  
8822:     if isinstance(obj, (np.integer, np.floating)):  
8823:         return obj.item()  
8824:     raise TypeError(f"Type {type(obj)} not serializable")  
8825:  
8826: def safe_utf8_truncate(text: str, max_bytes: int) -> str:  
8827:     """
```

```
8828:     Safely truncate text to max_bytes without cutting multi-byte UTF-8 characters.
8829:
8830:     Inputs:
8831:         text (str): Input text to truncate
8832:         max_bytes (int): Maximum number of UTF-8 bytes
8833:     Outputs:
8834:         str: Truncated text that is valid UTF-8
8835:     """
8836:     if not text:
8837:         return text
8838:     encoded = text.encode("utf-8")
8839:     if len(encoded) <= max_bytes:
8840:         return text
8841:     return encoded[:max_bytes].decode("utf-8", "ignore")
8842:
8843: def canonical_timestamp() -> str:
8844:     """
8845:     Generate ISO-8601 UTC timestamp with Z suffix for canonical timestamping.
8846:
8847:     Inputs:
8848:         None
8849:     Outputs:
8850:         str: ISO-8601 formatted UTC timestamp ending with 'Z'
8851:     """
8852:     return datetime.now(timezone.utc).isoformat().replace('+00:00', 'Z')
8853:
8854: def filter_empty_sentences(sentences: List[str]) -> List[str]:
8855:     """
8856:     Filter out empty or whitespace-only sentences.
8857:
8858:     Inputs:
8859:         sentences (List[str]): List of sentence strings
8860:     Outputs:
8861:         List[str]: Filtered list containing only non-empty sentences
8862:     """
8863:     return [s for s in sentences if s.strip()]
8864:
8865: # =====
8866: # CANONICAL ERROR CLASSES
8867: # =====
8868:
8869: class CanonicalError(Exception):
8870:     """Base class for canonical pipeline errors"""
8871:     pass
8872:
8873: class ValidationError(CanonicalError):
8874:     """Raised when input validation fails"""
8875:     pass
8876:
8877: class ProcessingError(CanonicalError):
8878:     """Raised when processing step fails"""
8879:     pass
8880:
8881: class SerializationError(CanonicalError):
8882:     """Raised when serialization fails"""
8883:     pass
```

```
8884:  
8885: # =====  
8886: # ENUMS Y TIPOS  
8887: # =====  
8888:  
8889: class ChunkType(Enum):  
8890:     DIAGNOSTICO = "diagnostico"  
8891:     ESTRATEGIA = "estrategia"  
8892:     METRICA = "metrica"  
8893:     FINANCIERO = "financiero"  
8894:     NORMATIVO = "normativo"  
8895:     OPERATIVO = "operativo"  
8896:     EVALUACION = "evaluacion"  
8897:     MIXTO = "mixto"  
8898:  
8899: class CausalRelationType(Enum):  
8900:     DIRECT_CAUSE = "direct_cause"  
8901:     INDIRECT_CAUSE = "indirect_cause"  
8902:     CONDITIONAL = "conditional"  
8903:     ENABLING = "enabling"  
8904:     PREVENTING = "preventing"  
8905:     CORRELATIONAL = "correlational"  
8906:     TEMPORAL_PRECEDENCE = "temporal_precedence"  
8907:  
8908: class PolicyEntityRole(Enum):  
8909:     EXECUTOR = "executor"  
8910:     BENEFICIARY = "beneficiary"  
8911:     REGULATOR = "regulator"  
8912:     FUNDER = "funder"  
8913:     STAKEHOLDER = "stakeholder"  
8914:     EVALUATOR = "evaluator" #  
8915:  
8916: # =====  
8917: # ESTRUCTURAS DE DATOS  
8918: # =====  
8919:  
8920: @dataclass  
8921: class CausalEvidence:  
8922:     dimension: str  
8923:     category: str  
8924:     matches: List[str]  
8925:     confidence: float  
8926:     context_span: Tuple[int, int]  
8927:     implicit_indicators: List[str]  
8928:     causal_type: CausalRelationType  
8929:     strength_score: float  
8930:     mechanisms: List[str] = field(default_factory=list)  
8931:     confounders: List[str] = field(default_factory=list)  
8932:     mediators: List[str] = field(default_factory=list)  
8933:     moderators: List[str] = field(default_factory=list)  
8934:  
8935: @dataclass  
8936: class PolicyEntity:  
8937:     entity_type: str  
8938:     text: str  
8939:     normalized_form: str
```

```
8940:     context_role: PolicyEntityRole
8941:     confidence: float
8942:     span: Tuple[int, int]
8943:     relationships: List[Tuple[str, str, float]] = field(default_factory=list)
8944:     attributes: Dict[str, Any] = field(default_factory=dict)
8945:     mentioned_count: int = 1
8946:
8947: @dataclass
8948: class Provenance:
8949:     """
8950:         Canonical provenance tracking for policy chunks.
8951:         Ensures 100% traceability to source document page and section.
8952:     """
8953:     page_number: int
8954:     section_header: Optional[str] = None
8955:     bbox: Optional[Tuple[float, float, float, float]] = None # (x0, y0, x1, y1)
8956:     span_in_page: Optional[Tuple[int, int]] = None
8957:     source_file: Optional[str] = None
8958:
8959: @dataclass
8960: class CrossDocumentReference:
8961:     target_section: str
8962:     reference_type: str
8963:     confidence: float
8964:     semantic_linkage: float
8965:     context_bridge: str
8966:     alignment_score: float = 0.0
8967:     bidirectional: bool = False
8968:     distance_in_doc: int = 0
8969:
8970: @dataclass
8971: class StrategicContext:
8972:     policy_intent: str
8973:     implementation_phase: str
8974:     geographic_scope: str
8975:     temporal_horizon: str
8976:     budget_linkage: str
8977:     risk_factors: List[str]
8978:     success_indicators: List[str]
8979:     alignment_with_sdg: List[str] = field(default_factory=list)
8980:     stakeholder_map: Dict[str, List[str]] = field(default_factory=dict)
8981:     policy_coherence_score: float = 0.0
8982:     intervention_logic_chain: List[str] = field(default_factory=list)
8983:
8984: @dataclass
8985: class ArgumentStructure:
8986:     claims: List[Tuple[str, float]]
8987:     evidence: List[Tuple[str, float]]
8988:     warrants: List[Tuple[str, float]]
8989:     backing: List[Tuple[str, float]]
8990:     rebuttals: List[Tuple[str, float]]
8991:     structure_type: str
8992:     strength_score: float
8993:     logical_coherence: float
8994:
8995: @dataclass
```

```
8996: class TemporalDynamics:
8997:     temporal_markers: List[Tuple[str, str, int]]
8998:     sequence_flow: List[Tuple[str, str, float]]
8999:     dependencies: List[Tuple[str, str, str]]
9000:     milestones: List[Dict[str, Any]]
9001:     temporal_coherence: float
9002:     causality_direction: str
9003:
9004: @dataclass
9005: class SmartPolicyChunk:
9006:     chunk_id: str
9007:     document_id: str
9008:     content_hash: str
9009:     text: str
9010:     normalized_text: str
9011:     semantic_density: float
9012:     section_hierarchy: List[str]
9013:     document_position: Tuple[int, int]
9014:     chunk_type: ChunkType
9015:     causal_chain: List[CausalEvidence]
9016:     policy_entities: List[PolicyEntity]
9017:     implicit_assumptions: List[Tuple[str, float]]
9018:     contextual_presuppositions: List[Tuple[str, float]]
9019:     policy_area_id: str # PA01-PA10 canonical code
9020:     dimension_id: str # DIM01-DIM06 canonical code
9021:
9022:     # SOTA Provenance Tracking (optional fields after required)
9023:     provenance: Optional[Provenance] = None
9024:     argument_structure: Optional[ArgumentStructure] = None
9025:     temporal_dynamics: Optional[TemporalDynamics] = None
9026:     discourse_markers: List[Tuple[str, str]] = field(default_factory=list)
9027:     rhetorical_patterns: List[str] = field(default_factory=list)
9028:
9029:     cross_references: List[CrossDocumentReference] = field(default_factory=list)
9030:     strategic_context: Optional[StrategicContext] = None
9031:     related_chunks: List[Tuple[str, float]] = field(default_factory=list)
9032:
9033:     confidence_metrics: Dict[str, float] = field(default_factory=dict)
9034:     coherence_score: float = 0.0
9035:     completeness_index: float = 0.0
9036:     strategic_importance: float = 0.0
9037:     information_density: float = 0.0
9038:     actionability_score: float = 0.0
9039:
9040:     semantic_embedding: Optional[np.ndarray] = None
9041:     policy_embedding: Optional[np.ndarray] = None
9042:     causal_embedding: Optional[np.ndarray] = None
9043:     temporal_embedding: Optional[np.ndarray] = None
9044:
9045:     knowledge_graph_nodes: List[str] = field(default_factory=list)
9046:     knowledge_graph_edges: List[Tuple[str, str, str, float]] = field(default_factory=list)
9047:
9048:     topic_distribution: Dict[str, float] = field(default_factory=dict)
9049:     key_phrases: List[Tuple[str, float]] = field(default_factory=list)
9050:
9051:     processing_timestamp: str = field(default_factory=canonical_timestamp)
```

```
9052:     pipeline_version: str = "SMART-CHUNK-3.0-FINAL"
9053:     extraction_methodology: str = "COMPREHENSIVE_STRATEGIC_ANALYSIS"
9054:     model_versions: Dict[str, str] = field(default_factory=dict) #
9055:
9056:     _CHUNK_ID_PATTERN = re.compile(r"^\w{1-9}\w{10}-\w{1-6}$")
9057:
9058:     def __post_init__(self) -> None:
9059:         """Validate canonical identifiers for irrigation compatibility."""
9060:         if not self.policy_area_id or not self.dimension_id:
9061:             raise ValueError("policy_area_id and dimension_id are required for SmartPolicyChunk")
9062:
9063:         match = self._CHUNK_ID_PATTERN.match(self.chunk_id)
9064:         if not match:
9065:             raise ValueError(
9066:                 f"Invalid chunk_id '{self.chunk_id}'. Expected format PA{{01-10}}-DIM{{01-06}}."
9067:             )
9068:
9069:         pa_code = f"PA{match.group(1)}"
9070:         dim_code = f"DIM{match.group(2)}"
9071:         if pa_code != self.policy_area_id:
9072:             raise ValueError(
9073:                 f"chunk_id {self.chunk_id} mismatches policy_area_id {self.policy_area_id}"
9074:             )
9075:         if dim_code != self.dimension_id:
9076:             raise ValueError(
9077:                 f"chunk_id {self.chunk_id} mismatches dimension_id {self.dimension_id}"
9078:             )
9079:
9080: # =====
9081: # CONFIGURACIÃN COMPLETA DEL SISTEMA
9082: # =====
9083:
9084: class SmartChunkConfig:
9085:     # ParÃ;metros de chunking calibrados
9086:     MIN_CHUNK_SIZE = 300
9087:     MAX_CHUNK_SIZE = 2000
9088:     OPTIMAL_CHUNK_SIZE = 800
9089:     OVERLAP_SIZE = 200
9090:
9091:     # Umbrales semÃ;nticos
9092:     SEMANTIC_COHERENCE_THRESHOLD = 0.72
9093:     CROSS_REFERENCE_MIN_SIMILARITY = 0.65
9094:     CAUSAL_CHAIN_MIN_CONFIDENCE = 0.60
9095:     ENTITY_EXTRACTION_THRESHOLD = 0.55
9096:
9097:     # ParÃ;metros de ventana de contexto
9098:     MIN_CONTEXT_WINDOW = 400
9099:     MAX_CONTEXT_WINDOW = 1200
9100:     CONTEXT_EXPANSION_FACTOR = 1.5
9101:
9102:     # Clustering y agrupaciÃ;n
9103:     DBSCAN_EPS = 0.25
9104:     DBSCAN_MIN_SAMPLES = 2
9105:     HIERARCHICAL_CLUSTER_THRESHOLD = 0.70
9106:
9107:     # AnÃ;lisis causal
```

```
9108:     CAUSAL_CHAIN_MAX_GAP = 3
9109:     TRANSITIVE_CLOSURE_DEPTH = 4
9110:     CAUSAL_MECHANISM_MIN_SUPPORT = 0.50
9111:
9112:     # TÃ³picos y temas
9113:     N_TOPICS_LDA = 15
9114:     MIN_TOPIC_PROBABILITY = 0.15
9115:
9116:     # MÃ©tricas de calidad
9117:     MIN_INFORMATION_DENSITY = 0.40
9118:     MIN_COHERENCE_SCORE = 0.55
9119:     MIN_COMPLETENESS_INDEX = 0.60
9120:     MIN_STRATEGIC_IMPORTANCE = 0.45
9121:
9122:     # DeduplicaciÃ³n
9123:     DEDUPLICATION_THRESHOLD = 0.88
9124:     NEAR_DUPLICATE_THRESHOLD = 0.92 #
9125:
9126: # =====
9127: # SISTEMAS AUXILIARES COMPLETOS
9128: # =====
9129:
9130: class ContextPreservationSystem:
9131:     """Sistema de preservaciÃ³n de contexto estratÃ©gico"""
9132:
9133:     def __init__(self, parent_system):
9134:         self.parent = parent_system
9135:         self.logger = logging.getLogger(self.__class__.__name__)
9136:
9137:     def preserve_strategic_context(
9138:         self,
9139:         text: str,
9140:         structural_analysis: Dict,
9141:         global_topics: Dict
9142:     ) -> List[Dict]:
9143:         """
9144:             CANONICAL_SOTA: Preserve strategic context using EmbeddingPolicyProducer.
9145:
9146:                 Derives breakpoints from canonical chunks with PDM structure awareness.
9147:                 Replaces internal breakpoint logic with SOTA semantic chunking.
9148:                 """
9149:                 # Use canonical chunker with doc_id/title from structural analysis
9150:                 chunks = self.parent._spc_embed.process_document(
9151:                     text,
9152:                     {
9153:                         "doc_id": structural_analysis.get("doc_id", "unknown"),
9154:                         "title": structural_analysis.get("title", "N/A")
9155:                     }
9156:                 )
9157:
9158:                 # Build segments from canonical chunks (position is ordinal; approximate char spans)
9159:                 segments = []
9160:                 offset = 0
9161:                 for ch in chunks:
9162:                     ch_text = self.parent._spc_embed.get_chunk_text(ch)
9163:                     start = offset
```

```

9164:         end = start + len(ch_text)
9165:         offset = end
9166:
9167:         segment = {
9168:             "text": ch_text,
9169:             "context": ch_text, # Can expand context if needed
9170:             "position": (start, end),
9171:             "context_window": (start, end),
9172:             "semantic_coherence": 0.0, # Filled by coherence calculation if needed
9173:             "topic_alignment": self._calculate_topic_alignment(ch_text, global_topics),
9174:             "pdq_context": self.parent._spc_embed.get_chunk_pdq_context(ch),
9175:             "metadata": self.parent._spc_embed.get_chunk_metadata(ch),
9176:         }
9177:         segments.append(segment)
9178:
9179:     return segments
9180:
9181: def _identify_semantic_breakpoints(self, text: str, structural_analysis: Dict) -> List[int]:
9182:     """Identificar puntos de ruptura semántica"""
9183:     breakpoints = [0]
9184:
9185:     # Usar lÃ¡mites de secciÃ³n
9186:     for section in structural_analysis.get('section_hierarchy', []):
9187:         if 'line_number' in section:
9188:             pos = self._line_to_position(text, section['line_number'])
9189:             breakpoints.append(pos)
9190:
9191:     # Usar puntos de quiebre estratÃ©gico
9192:     for bp in structural_analysis.get('strategic_breakpoints', []):
9193:         breakpoints.append(bp['position'])
9194:
9195:     # Agregar lÃ¡mites de pÃ¡rrafo significativos
9196:     paragraphs = text.split('\n\n')
9197:     current_pos = 0
9198:     for para in paragraphs:
9199:         if len(para) > self.parent.config.MIN_CHUNK_SIZE:
9200:             breakpoints.append(current_pos)
9201:             current_pos += len(para) + 2
9202:
9203:     breakpoints.append(len(text))
9204:     return sorted(list(set(breakpoints)))
9205:
9206: def _line_to_position(self, text: str, line_number: int) -> int:
9207:     """Convertir nÃºmero de lÃnea a posiciÃ³n en texto"""
9208:     lines = text.split('\n')
9209:     position = 0
9210:     for i in range(min(line_number, len(lines))):
9211:         position += len(lines[i]) + 1
9212:     return position
9213:
9214: def _calculate_segment_coherence(self, segment_text: str) -> float:
9215:     """
9216:         CANONICAL SOTA: Calculate coherence using batch embeddings.
9217:
9218:         Replaces per-sentence embedding calls with efficient batching.
9219:         No per-sentence model churn.

```

```
9220:  
9221:     Inputs:  
9222:         segment_text (str): Text segment to analyze  
9223:     Outputs:  
9224:         float: Coherence score between 0.0 and 1.0  
9225:         """  
9226:     if len(segment_text) < 50:  
9227:         return 0.0  
9228:  
9229:     sentences = filter_empty_sentences(re.split(r'[.!?]+', segment_text))  
9230:     if len(sentences) < 2:  
9231:         return 0.5  
9232:  
9233:     # CANONICAL SOTA: Batch embeddings for efficiency  
9234:     embs = self.parent._generate_embeddings_for_corpus(sentences, batch_size=64)  
9235:  
9236:     # Pairwise cosine similarity between consecutive sentences  
9237:     sims = np.sum(embs[:-1] * embs[1:], axis=1) / (  
9238:         np.linalg.norm(embs[:-1], axis=1) * np.linalg.norm(embs[1:], axis=1) + 1e-8  
9239:     )  
9240:  
9241:     return float(np.mean(sims)) if sims.size else 0.5  
9242:  
9243: def _calculate_topic_alignment(self, segment_text: str, global_topics: Dict) -> float:  
9244:     """Calcular alineaciÃ³n con tÃ³picos globales"""  
9245:     if not global_topics.get('keywords'):  
9246:         return 0.5  
9247:  
9248:     segment_lower = segment_text.lower()  
9249:     keyword_matches = 0  
9250:  
9251:     for keyword, _ in global_topics['keywords'][:20]:  
9252:         if keyword.lower() in segment_lower:  
9253:             keyword_matches += 1  
9254:  
9255:     return min(keyword_matches / 10.0, 1.0) #  
9256:  
9257:  
9258: class CausalChainAnalyzer:  
9259:     """Analizador de cadenas causales"""  
9260:  
9261:     def __init__(self, parent_system):  
9262:         self.parent = parent_system  
9263:         self.logger = logging.getLogger(self.__class__.__name__)  
9264:  
9265:     def extract_complete_causal_chains(  
9266:         self,  
9267:         segments: List[Dict],  
9268:         knowledge_graph: Dict  
9269:     ) -> List[Dict]:  
9270:         """Extraer cadenas causales completas"""  
9271:         causal_chains = []  
9272:  
9273:         for segment in segments:  
9274:             chains = self._extract_segment_causal_chains(segment, knowledge_graph)  
9275:             causal_chains.extend(chains)
```

```

9276:
9277:     # Conectar cadenas entre segmentos
9278:     connected_chains = self._connect_cross_segment_chains(causal_chains)
9279:
9280:     return connected_chains
9281:
9282: def _extract_segment_causal_chains(self, segment: Dict, kg: Dict) -> List[Dict]:
9283:     """Extraer cadenas causales de un segmento"""
9284:     text = segment['text']
9285:     chains = []
9286:
9287:     # Patrones causales complejos
9288:     causal_patterns = [
9289:         r'si\s+([^\s]+)\s+entonces\s+([^\s]+)', 'conditional'),
9290:         r'debido\s+a\s+([^\s]+)\s+([^\s]+)', 'direct_cause'),
9291:         r'([^\s]+)\s+permite\s+([^\s]+)', 'enabling'),
9292:         r'([^\s]+)\s+genera\s+([^\s]+)', 'generation'),
9293:         r'para\s+([^\s]+)\s+se\s+requiere\s+([^\s]+)', 'requirement'),
9294:         r'([^\s]+)\s+resulta\s+en\s+([^\s]+)', 'result')
9295:     ]
9296:
9297:     for pattern, chain_type in causal_patterns:
9298:         matches = re.finditer(pattern, text, re.IGNORECASE)
9299:         for match in matches:
9300:             chain = {
9301:                 'type': chain_type,
9302:                 'antecedent': match.group(1).strip(),
9303:                 'consequent': match.group(2).strip() if match.lastindex >= 2 else '',
9304:                 'position': match.span(),
9305:                 'segment_id': id(segment),
9306:                 'confidence': self._calculate_causal_confidence(match.group(0), text)
9307:             }
9308:             chains.append(chain)
9309:
9310:     return chains
9311:
9312: def _connect_cross_segment_chains(self, chains: List[Dict]) -> List[Dict]:
9313:     """Conectar cadenas causales entre segmentos"""
9314:     if len(chains) < 2:
9315:         return chains
9316:
9317:     # Construir grafo de cadenas
9318:     G = nx.DiGraph()
9319:
9320:     for i, chain in enumerate(chains):
9321:         G.add_node(i, **chain)
9322:
9323:     # Conectar cadenas relacionadas
9324:     for i in range(len(chains)):
9325:         for j in range(i + 1, len(chains)):
9326:             similarity = self._calculate_chain_similarity(chains[i], chains[j])
9327:             if similarity > 0.7:
9328:                 G.add_edge(i, j, weight=similarity)
9329:
9330:     # Enriquecer cadenas con conexiones
9331:     for i, chain in enumerate(chains):

```

```

9332:         chain['connections'] = list(G.neighbors(i))
9333:         chain['centrality'] = nx.degree_centrality(G).get(i, 0)
9334:
9335:     return chains
9336:
9337: def _calculate_causal_confidence(self, match_text: str, context: str) -> float:
9338:     """Calcular confianza de relaciÃ³n causal"""
9339:     confidence = 0.5
9340:
9341:     # Indicadores de confianza alta
9342:     high_confidence_terms = ['garantiza', 'asegura', 'determina', 'causa directamente']
9343:     for term in high_confidence_terms:
9344:         if term in match_text.lower() or term in context.lower():
9345:             confidence = max(confidence, 0.85)
9346:
9347:     # Indicadores de confianza media
9348:     medium_confidence_terms = ['permite', 'facilita', 'contribuye', 'apoya']
9349:     for term in medium_confidence_terms:
9350:         if term in match_text.lower():
9351:             confidence = max(confidence, 0.65)
9352:
9353:     # Indicadores de incertidumbre
9354:     uncertainty_terms = ['puede', 'podrÃ¡-a', 'posiblemente', 'eventualmente']
9355:     for term in uncertainty_terms:
9356:         if term in match_text.lower():
9357:             confidence = min(confidence, 0.45)
9358:
9359:     return confidence
9360:
9361: def _calculate_chain_similarity(self, chain1: Dict, chain2: Dict) -> float:
9362:     """
9363:         CANONICAL SOTA: Calculate chain similarity via batch embeddings.
9364:
9365:         Replaces individual embedding calls with efficient batching.
9366:
9367:         Inputs:
9368:             chain1 (Dict): First causal chain
9369:             chain2 (Dict): Second causal chain
9370:         Outputs:
9371:             float: Similarity score between 0.0 and 1.0
9372:     """
9373:     # Build text representations of chains
9374:     texts = [
9375:         f"{chain1.get('antecedent', '')} {chain1.get('consequent', '')}",
9376:         f"{chain2.get('antecedent', '')} {chain2.get('consequent', '')}",
9377:     ]
9378:
9379:     # CANONICAL SOTA: Batch embeddings for efficiency
9380:     embs = self.parent._generate_embeddings_for_corpus(texts, batch_size=2)
9381:
9382:     # Cosine similarity
9383:     sim = float(np.dot(embs[0], embs[1]) / (
9384:         np.linalg.norm(embs[0]) * np.linalg.norm(embs[1]) + 1e-8
9385:     ))
9386:
9387:     # PenalizaciÃ³n por tipo de relaciÃ³n diferente

```

```

9388:     if chain1.get('type') != chain2.get('type'):
9389:         sim *= 0.9
9390:
9391:     return sim
9392:
9393:
9394: class KnowledgeGraphBuilder:
9395:     """Constructor de grafo de conocimiento para polÃtica pÃblica"""
9396:
9397:     def __init__(self, parent_system):
9398:         self.parent = parent_system
9399:         self.logger = logging.getLogger(self.__class__.__name__)
9400:
9401:     def build_policy_knowledge_graph(self, text: str) -> Dict[str, Any]:
9402:         """Construir grafo de conocimiento de polÃtica pÃblica"""
9403:         G = nx.DiGraph()
9404:         entities = self._extract_all_entities(text)
9405:         relations = self._extract_all_relations(text)
9406:         concepts = self._extract_key_concepts(text)
9407:
9408:         # AÃ±adir entidades como nodos
9409:         for entity in entities:
9410:             G.add_node(entity['normalized_form'], type=entity['entity_type'], role=entity['context_role'].value, confidence=entity['confidence'])
9411:
9412:         # AÃ±adir relaciones como aristas
9413:         for relation in relations:
9414:             source = self.parent._normalize_entity(relation['source'])
9415:             target = self.parent._normalize_entity(relation['target'])
9416:             if source in G.nodes and target in G.nodes:
9417:                 G.add_edge(source, target, type=relation['type'], confidence=relation['confidence'])
9418:
9419:         # MÃ©tricas del grafo
9420:         metrics = {
9421:             'num_nodes': G.number_of_nodes(),
9422:             'num_edges': G.number_of_edges(),
9423:             'density': nx.density(G) if G.number_of_nodes() > 0 else 0,
9424:             'components': nx.number_weakly_connected_components(G) if G.number_of_nodes() > 0 else 0
9425:         }
9426:
9427:         return {
9428:             'graph': G,
9429:             'entities': entities,
9430:             'relations': relations,
9431:             'concepts': concepts,
9432:             'metrics': metrics
9433:         }
9434:
9435:     def _extract_all_entities(self, text: str) -> List[Dict]:
9436:         """Extraer todas las entidades del texto"""
9437:         entities = []
9438:
9439:         # Patrones de entidades por tipo
9440:         entity_patterns = {
9441:             'organization': [
9442:                 r'(?:(Ministerio|SecretarÃ-a|Departamento|DirecciÃ³n|Instituto)\s+(:de|del?))\s+[A-ZÃ\u0121\u211A\u215A\u223A\u232A\u2211][a-zÃ\u0121Ã\u0131Ã\u0151\s]+',
9443:                 r'(?:(AlcaldÃ-a|GobernaciÃ³n|Prefectura)\s+(:de|del?))\s+[A-ZÃ\u0121\u211A\u215A\u223A\u232A\u2211][a-zÃ\u0121Ã\u0131Ã\u0151\s]+'
9444:             ]
9445:         }

```

```

9444:
9445:     ],
9446:     'program': [
9447:         r'(?:Programa|Plan|Proyecto)\s+(?:de|del?|para)\s+[A-ZÀ\201Ã\211Ã\215Ã\223Ã\232Ã\221] [a-zA-À-Ã³Ã°Ã±\s]+',
9448:     ],
9449:     'legal_framework': [
9450:         r'(?:Ley|Decreto|ResoluciÃ³n|Acuerdo)\s+No?\s+[\d]+(?: de \d{4})?',
9451:     ]
9452:
9453:     for entity_type, patterns in entity_patterns.items():
9454:         for pattern in patterns:
9455:             matches = re.findall(pattern, text)
9456:             for match in matches:
9457:                 role = self.parent._infer_entity_role(match.group(0), text)
9458:                 entities.append({
9459:                     'text': match.group(0),
9460:                     'normalized_form': self.parent._normalize_entity(match.group(0)),
9461:                     'entity_type': entity_type,
9462:                     'context_role': role,
9463:                     'confidence': 0.8,
9464:                     'span': match.span()
9465:                 })
9466:
9467:     return entities
9468:
9469: def _extract_all_relations(self, text: str) -> List[Dict]:
9470:     """Extraer todas las relaciones del texto"""
9471:     relations = []
9472:
9473:     # Patrones de relaciones
9474:     relation_patterns = [
9475:         (r'([^\s]+)\s+es\s+responsable\s+de\s+([^\s]+)', 'responsible_for', 0.9),
9476:         (r'([^\s]+)\s+ejecutarÃ;|\s+([^\s]+)', 'executes', 0.85),
9477:         (r'([^\s]+)\s+beneficiarÃ;|\s+a\s+([^\s]+)', 'benefits', 0.8),
9478:         (r'([^\s]+)\s+financiarÃ;|\s+([^\s]+)', 'funds', 0.85)
9479:     ]
9480:
9481:     for pattern, rel_type, confidence in relation_patterns:
9482:         matches = re.findall(pattern, text, re.IGNORECASE)
9483:         for match in matches:
9484:             if match.lastindex >= 2:
9485:                 relations.append({
9486:                     'source': match.group(1).strip(),
9487:                     'target': match.group(2).strip(),
9488:                     'type': rel_type,
9489:                     'confidence': confidence,
9490:                     'context': match.group(0)
9491:                 })
9492:
9493:     return relations
9494:
9495: def _extract_key_concepts(self, text: str) -> List[str]:
9496:     """
9497:         Extract key concepts from noun chunks and key phrases.
9498:
9499:         Inputs:

```

```

9500:         text (str): Input text to analyze
9501:     Outputs:
9502:         List[str]: List of key concepts (max 30)
9503:     """
9504:     concepts = []
9505:     if self.parent.nlp:
9506:         # Safe UTF-8 truncation to avoid cutting multi-byte characters
9507:         truncated_text = safe_utf8_truncate(text, 200000)
9508:         doc = self.parent.nlp(truncated_text)
9509:         concepts.extend([chunk.text for chunk in doc.noun_chunks][:50])
9510:
9511:     # Normalizar conceptos
9512:     concepts = list(set([c.lower().strip() for c in concepts]))
9513:     return concepts[:30]
9514:
9515:
9516: class TopicModeler:
9517:     """Modelador de tÃ³picos y temas"""
9518:
9519:     def __init__(self, parent_system):
9520:         self.parent = parent_system
9521:         self.logger = logging.getLogger(self.__class__.__name__)
9522:         self.tfidf_vectorizer = TfidfVectorizer(stop_words=self.parent._get_stopwords(), ngram_range=(1, 2), max_df=0.85, min_df=2)
9523:         self.lda_model = LatentDirichletAllocation(n_components=self.parent.config.N_TOPICS_LDA, random_state=42)
9524:
9525:     def _get_stopwords(self) -> List[str]:
9526:         """Obtener lista de stopwords en espaÃ±ol (placeholder)"""
9527:         # Una lista de stopwords mÃ¡s completa se usarÃ¡ en producciÃ³n
9528:         return ['el', 'la', 'los', 'las', 'un', 'una', 'unos', 'unas', 'y', 'o', 'de', 'a', 'en', 'por', 'con', 'para', 'del', 'al', 'que', 'se', 'es', 'son',
9529:         'han', 'como', 'mÃ¡s', 'pero', 'no', 'su', 'sus']
9530:
9531:     def extract_global_topics(self, text_list: List[str]) -> Dict[str, Any]:
9532:         """Extraer tÃ³picos globales mediante LDA"""
9533:         try:
9534:             # 1. Vectorizar
9535:             tfidf_matrix = self.tfidf_vectorizer.fit_transform(text_list)
9536:
9537:             # 2. Aplicar LDA
9538:             self.lda_model.fit(tfidf_matrix)
9539:             lda_output = self.lda_model.transform(tfidf_matrix)
9540:
9541:             # 3. Extraer tÃ³picos y palabras clave
9542:             feature_names = self.tfidf_vectorizer.get_feature_names_out()
9543:             topics = []
9544:
9545:             for topic_idx, topic in enumerate(self.lda_model.components_):
9546:                 top_features_ind = topic.argsort()[:-10 - 1:-1]
9547:                 top_features = [(feature_names[i], topic[i]) for i in top_features_ind]
9548:
9549:                 topics.append({
9550:                     'topic_id': topic_idx,
9551:                     'keywords': top_features,
9552:                     'weight': float(topic.sum())
9553:                 })
9554:
9555:             # Palabras clave globales

```

```
9555:     global_keywords = []
9556:     for topic in topics:
9557:         global_keywords.extend([kw[0] for kw in topic['keywords']])
9558:     keyword_counts = Counter(global_keywords)
9559:     top_keywords = keyword_counts.most_common(30)
9560:
9561:     return {
9562:         'topics': topics,
9563:         'keywords': top_keywords,
9564:         'topic_distribution': lda_output.mean(axis=0).tolist()
9565:     }
9566: except Exception as e:
9567:     self.logger.error(f"Error en extracciÃ³n de tÃ³picos: {e}")
9568: return {'topics': [], 'keywords': []} #
9569:
9570:
9571: class ArgumentAnalyzer:
9572:     """Analizador de estructura argumentativa (Toulmin)"""
9573:
9574:     def __init__(self, parent_system):
9575:         self.parent = parent_system
9576:         self.logger = logging.getLogger(self.__class__.__name__)
9577:
9578:     def analyze_arguments(self, causal_chains: List[Dict]) -> Dict[int, ArgumentStructure]:
9579:         """Analizar argumentos completos"""
9580:         argument_structures = {}
9581:
9582:         # Agrupar cadenas para formar argumentos
9583:         for idx, chain_group in enumerate(self._group_chains_by_proximity(causal_chains)):
9584:             structure = self._extract_argument_structure(chain_group)
9585:             if structure:
9586:                 argument_structures[idx] = structure
9587:
9588:         return argument_structures
9589:
9590:     def _group_chains_by_proximity(self, chains: List[Dict], max_gap: int = 500) -> List[List[Dict]]:
9591:         """Agrupar cadenas causales por proximidad en el texto"""
9592:         if not chains:
9593:             return []
9594:
9595:         chains.sort(key=lambda x: x['position'][0])
9596:         groups = []
9597:         current_group = [chains[0]]
9598:
9599:         for i in range(1, len(chains)):
9600:             prev_end = chains[i-1]['position'][1]
9601:             current_start = chains[i]['position'][0]
9602:
9603:             if current_start - prev_end < max_gap:
9604:                 current_group.append(chains[i])
9605:             else:
9606:                 groups.append(current_group)
9607:                 current_group = [chains[i]]
9608:
9609:             if current_group:
9610:                 groups.append(current_group)
```

```
9611:         return groups
9612:
9613:
9614:     def _extract_argument_structure(self, chain_group: List[Dict]) -> Optional[ArgumentStructure]:
9615:         """Extraer los componentes del argumento (Claims, Evidence, Warrants)"""
9616:         if not chain_group:
9617:             return None
9618:
9619:         claims = []
9620:         evidence = []
9621:         warrants = []
9622:
9623:         full_text = ' '.join([f'{c.get("antecedent", "")} {c.get("consequent", "")}' for c in chain_group])
9624:
9625:         # Claims: Resultados (consequents) o afirmaciones directas
9626:         for chain in chain_group:
9627:             if chain.get('type') in ['result', 'generation']:
9628:                 claims.append((chain.get('consequent', ''), chain.get('confidence', 0.5)))
9629:
9630:         # Evidence: Antecedentes o referencias a datos/normas
9631:         for chain in chain_group:
9632:             if chain.get('type') in ['direct_cause', 'conditional', 'requirement']:
9633:                 evidence.append((chain.get('antecedent', ''), chain.get('confidence', 0.5)))
9634:
9635:         # Warrants: Conexiones causales implÃ¡-citas o explÃ¡-citas de alta confianza
9636:         warrants.extend(self._identify_warrants(chain_group))
9637:
9638:         # Backing and Rebuttals (Simplificado: Requiere modelo avanzado)
9639:         backing = []
9640:         rebuttals = []
9641:
9642:         # Estructura y Fuerza
9643:         structure_type = self._determine_structure_type(claims, evidence)
9644:         strength_score = self._calculate_argument_strength(claims, evidence, warrants)
9645:         logical_coherence = self._assess_logical_coherence(claims, evidence)
9646:
9647:         return ArgumentStructure(
9648:             claims=claims[:5],
9649:             evidence=evidence[:5],
9650:             warrants=warrants[:3],
9651:             backing=backing,
9652:             rebuttals=rebuttals,
9653:             structure_type=structure_type,
9654:             strength_score=strength_score,
9655:             logical_coherence=logical_coherence
9656:         )
9657:
9658:     def _identify_warrants(self, chain_group: List[Dict]) -> List[Tuple[str, float]]:
9659:         """Identificar warrants (garantÃ‐as/conexiones)"""
9660:         warrants = []
9661:         for chain in chain_group:
9662:             if chain.get('confidence', 0.5) > 0.7:
9663:                 warrants.append((f'ConexiÃ³n causal tipo: {chain.get("type")}', chain.get('confidence', 0.5)))
9664:         return warrants
9665:
9666:     def _determine_structure_type(self, claims: List, evidence: List) -> str:
```

```
9667:     """Determinar el tipo de estructura argumentativa"""
9668:     if len(claims) > 1 and len(evidence) >= 1:
9669:         return 'multiple_claims_supported'
9670:     elif len(claims) == 1 and len(evidence) >= 1:
9671:         return 'simple_supported'
9672:     elif not evidence and claims:
9673:         return 'assertion_only'
9674:     elif len(claims) >= 1 and any(c.lower().startswith(('segÃ³n', 'de acuerdo con')) for c, _ in evidence):
9675:         return 'evidence_based'
9676:     else:
9677:         return 'balanced'
9678:
9679: def _calculate_argument_strength(self, claims: List, evidence: List, warrants: List) -> float:
9680:     """Calcular fuerza del argumento"""
9681:     if not claims:
9682:         return 0.0
9683:
9684:     claim_strength = np.mean([conf for _, conf in claims]) if claims else 0
9685:     evidence_strength = np.mean([conf for _, conf in evidence]) if evidence else 0
9686:     warrant_strength = np.mean([conf for _, conf in warrants]) if warrants else 0
9687:
9688:     # PonderaciÃ³n: evidencia mÃ¡s importante que claims
9689:     strength = (claim_strength * 0.3 + evidence_strength * 0.5 + warrant_strength * 0.2)
9690:
9691:     # Penalizar argumentos sin evidencia
9692:     if not evidence:
9693:         strength *= 0.5
9694:
9695:     return min(strength, 1.0)
9696:
9697: def _assess_logical_coherence(self, claims: List, evidence: List) -> float:
9698:     """Evaluar coherencia lÃ³gica del argumento"""
9699:     if not claims or not evidence:
9700:         return 0.0
9701:
9702:     # Coherencia basada en similitud semÃ¡ntica entre claims y evidencia
9703:     all_texts = [c for c, _ in claims] + [e for e, _ in evidence]
9704:     if len(all_texts) < 2:
9705:         return 0.5
9706:
9707:     # Use batch embedding for efficiency
9708:     embeddings = self.parent._generate_embeddings_for_corpus(all_texts, batch_size=64)
9709:
9710:     # Vectorized cosine similarity computation (no sklearn dependency)
9711:     # Normalize embeddings for efficient dot product = cosine similarity
9712:     norms = np.linalg.norm(embeddings, axis=1, keepdims=True)
9713:     normalized_embs = embeddings / (norms + 1e-8)
9714:     sim_matrix = np.dot(normalized_embs, normalized_embs.T)
9715:
9716:     # Tomar la similitud media (excluyendo la diagonal)
9717:     coherence = (np.sum(sim_matrix) - np.trace(sim_matrix)) / (len(sim_matrix)**2 - len(sim_matrix))
9718:
9719:     return min(max(coherence, 0.0), 1.0)
9720:
9721:
9722: class TemporalAnalyzer:
```

```
9723: """Analizador de dinÃ¡mica temporal y secuencial"""
9724:
9725: def __init__(self, parent_system):
9726:     self.parent = parent_system
9727:     self.logger = logging.getLogger(self.__class__.__name__)
9728:
9729: def analyze_temporal_dynamics(self, causal_chains: List[Dict]) -> Dict[int, Optional[TemporalDynamics]]:
9730:     """Analizar dinÃ¡mica temporal completa"""
9731:     temporal_structures = {}
9732:
9733:     # Agrupar cadenas para anÃ;lisis temporal
9734:     for idx, chain_group in enumerate(self.parent.argument_analyzer._group_chains_by_proximity(causal_chains)):
9735:         structure = self._extract_temporal_structure(chain_group)
9736:         if structure:
9737:             temporal_structures[idx] = structure
9738:
9739:     return temporal_structures
9740:
9741: def _extract_temporal_structure(self, chain_group: List[Dict]) -> Optional[TemporalDynamics]:
9742:     """Extraer marcadores, secuencias y dependencias temporales"""
9743:     if not chain_group:
9744:         return None
9745:
9746:     temporal_markers = []
9747:     sequence_flow = []
9748:     dependencies = []
9749:     milestones = []
9750:
9751:     # Extraer marcadores de tiempo
9752:     for chain in chain_group:
9753:         # Reutilizar el analizador temporal de la clase principal
9754:         text_context = f'{chain.get("antecedent", "")} {chain.get("consequent", "")}'
9755:         temp_info = self.parent._analyze_temporal_structure(text_context)
9756:
9757:         for marker in temp_info.get('time_markers', []):
9758:             temporal_markers.append((marker['text'], marker['type'], marker['position'][0]))
9759:
9760:         for seq in temp_info.get('sequences', []):
9761:             sequence_flow.append((seq['marker'], str(seq['order']), 0.8))
9762:
9763:     # Extraer dependencias causales con implicaciÃ³n temporal
9764:     for chain in chain_group:
9765:         if chain.get('type') in ['conditional', 'direct_cause', 'requirement']:
9766:             dependencies.append((
9767:                 chain.get('antecedent', ''),
9768:                 chain.get('consequent', ''),
9769:                 'prerequisite' if chain.get('type') == 'requirement' else 'temporal_precedence'
9770:             ))
9771:         elif chain.get('type') == 'conditional':
9772:             dependencies.append((
9773:                 chain.get('antecedent', ''),
9774:                 chain.get('consequent', ''),
9775:                 'conditional'
9776:             ))
9777:
9778:     # Extraer hitos
```

```

9779:     milestone_patterns = [
9780:         r'meta\s+(?:de|para)\s+([^.]+)',
9781:         r'lograr\s+([^.]+)\s+(?:en|para)\s+(20\d{2})',
9782:         r'alcanzar\s+([^.]+)'
9783:     ]
9784:
9785:     for chain in chain_group:
9786:         text = f"{chain.get('antecedent', '')} {chain.get('consequent', '')}"
9787:         for pattern in milestone_patterns:
9788:             matches = re.findall(pattern, text, re.IGNORECASE)
9789:             for match in matches[:2]:
9790:                 milestones.append({
9791:                     'description': match if isinstance(match, str) else match[0],
9792:                     'target_date': match[1] if isinstance(match, tuple) and len(match) > 1 else None,
9793:                     'confidence': chain.get('confidence', 0.5)
9794:                 })
9795:
9796:     if not (temporal_markers or sequence_flow or dependencies or milestones):
9797:         return None
9798:
9799:     temporal_coherence = self._calculate_temporal_coherence(temporal_markers, sequence_flow)
9800:     causality_direction = self._determine_causality_direction(dependencies)
9801:
9802:     return TemporalDynamics(
9803:         temporal_markers=temporal_markers[:10],
9804:         sequence_flow=sequence_flow[:5],
9805:         dependencies=dependencies[:10],
9806:         milestones=milestones[:5],
9807:         temporal_coherence=temporal_coherence,
9808:         causality_direction=causality_direction
9809:     )
9810:
9811:     def _calculate_temporal_coherence(self, markers: List, flow: List) -> float:
9812:         """Calcular la coherencia de los marcadores temporales"""
9813:         # Simple mÃ¡trica basada en la presencia de orden y hitos
9814:         score = 0.0
9815:         if flow:
9816:             score += 0.5
9817:             if any(m[1] in ['year', 'month_year', 'period'] for m in markers):
9818:                 score += 0.5
9819:             return min(score, 1.0)
9820:
9821:     def _determine_causality_direction(self, dependencies: List[Tuple]) -> str:
9822:         """Determinar la direcciÃ³n dominante de la causalidad (forward/backward)"""
9823:         forward = sum(1 for _, _, t in dependencies if t == 'temporal_precedence')
9824:         backward = sum(1 for _, _, t in dependencies if t == 'backward')
9825:
9826:         if forward > backward:
9827:             return 'forward'
9828:         elif backward > forward:
9829:             return 'backward'
9830:         return 'mixed' #
9831:
9832:
9833: class DiscourseAnalyzer:
9834:     """Analizador de discurso"""

```

```
9835:
9836:     def __init__(self, parent_system):
9837:         self.parent = parent_system
9838:         self.logger = logging.getLogger(self.__class__.__name__)
9839:
9840:     def analyze_discourse(self, causal_chains: List[Dict]) -> Dict[int, Dict]:
9841:         """Analizar estructuras discursivas"""
9842:         discourse_structures = {}
9843:
9844:         for idx, chain_group in enumerate(self.parent.argument_analyzer._group_chains_by_proximity(causal_chains, max_gap=300)):
9845:             structure = self._extract_discourse_structure(chain_group)
9846:             if structure:
9847:                 discourse_structures[idx] = structure
9848:
9849:         return discourse_structures
9850:
9851:     def _group_chains_discursively(self, chains: List[Dict]) -> List[List[Dict]]:
9852:         """Agrupar cadenas por coherencia discursiva (reutiliza lógica de ArgumentAnalyzer)"""
9853:         return self.parent.argument_analyzer._group_chains_by_proximity(chains, max_gap=300)
9854:
9855:     def _extract_discourse_structure(self, chain_group: List[Dict]) -> Optional[Dict]:
9856:         """Extraer estructura discursiva del grupo"""
9857:         if not chain_group:
9858:             return None
9859:
9860:         full_text = ' '.join([f'{c.get("antecedent", "")} {c.get("consequent", "")}' for c in chain_group])
9861:
9862:         # Análisis de relaciones
9863:         relations = self.parent._extract_coherence_relations(full_text)
9864:
9865:         # Análisis retórico
9866:         rhetorical = self.parent._analyze_rhetorical_structure(full_text)
9867:
9868:         # Análisis de flujo de información
9869:         info_flow = self.parent._analyze_information_flow(full_text)
9870:
9871:         return {
9872:             'coherence_relations': relations[:10],
9873:             'rhetorical_moves': list(rhetorical.keys()),
9874:             'flow_metrics': info_flow,
9875:             'complexity_score': self._calculate_discourse_complexity(relations, rhetorical)
9876:         }
9877:
9878:     def _calculate_discourse_complexity(self, relations: List[Dict], rhetorical: Dict) -> float:
9879:         """Calcular complejidad discursiva"""
9880:         moves = []
9881:         for v in rhetorical.values():
9882:             moves.extend(v)
9883:
9884:         move_diversity = len(set(moves)) if moves else 0
9885:         relation_diversity = len(set(r['type'] for r in relations)) if relations else 0
9886:
9887:         complexity = (move_diversity / 5.0) * 0.5 + (relation_diversity / 8.0) * 0.5
9888:         return min(complexity, 1.0) #
```

```

9891: class StrategicIntegrator:
9892:     """Integrador de análisis multi-escala"""
9893:
9894:     def __init__(self, parent_system):
9895:         self.parent = parent_system
9896:         self.logger = logging.getLogger(self.__class__.__name__)
9897:
9898:     def integrate_strategic_units(
9899:         self,
9900:         causal_chains: List[Dict],
9901:         structural_analysis: Dict,
9902:         argument_structures: Dict,
9903:         temporal_structures: Dict,
9904:         discourse_structures: Dict,
9905:         global_topics: Dict,
9906:         global_kg: Dict
9907:     ) -> List[Dict]:
9908:         """Integrar todos los análisis en unidades estratégicas"""
9909:
9910:         # 1. Agrupar cadenas en unidades base (reutilizando la agrupación)
9911:         grouped_chains = self.parent.argument_analyzer._group_chains_by_proximity(causal_chains, max_gap=100)
9912:
9913:         strategic_units = []
9914:
9915:         for idx, chain_group in enumerate(grouped_chains):
9916:             full_text = ' '.join([f'{c.get("antecedent", '')} {c.get("consequent", '')}' for c in chain_group])
9917:             # Derive hierarchy for this segment
9918:             start_pos = chain_group[0]['position'][0] if chain_group and 'position' in chain_group[0] else 0
9919:             hierarchy = self._derive_hierarchy_for_segment(start_pos, structural_analysis)
9920:             # Unidades de integración
9921:             unit = {
9922:                 'index': idx,
9923:                 'text': full_text,
9924:                 'position': (chain_group[0]['position'][0], chain_group[-1]['position'][1]),
9925:                 'chains': chain_group,
9926:                 'argument_structure': argument_structures.get(idx),
9927:                 'temporal_dynamics': temporal_structures.get(idx),
9928:                 'discourse_structure': discourse_structures.get(idx),
9929:                 'hierarchy': hierarchy,
9930:                 'confidence': np.mean([c.get('confidence', 0.5) for c in chain_group]),
9931:                 'semantic_coherence': self.parent.context_preserver._calculate_segment_coherence(full_text)
9932:             }
9933:
9934:             strategic_units.append(unit)
9935:
9936:         # 2. Enriquecer con metadatos estratégicos
9937:         enriched_units = self._enrich_strategic_metadata(strategic_units)
9938:
9939:         # 3. Refinar lánquimes de las unidades (placeholder para refinamiento avanzado)
9940:
9941:         return enriched_units
9942:
9943:     def _derive_hierarchy_for_segment(self, start_pos: int, structural_analysis: Dict) -> List[str]:
9944:         """Derivar la jerarquía de sección para una posición en el texto"""
9945:         hierarchy = []
9946:         best_match = None

```

```

9947:     min_distance = float('inf')
9948:
9949:     for section in structural_analysis.get('section_hierarchy', []):
9950:         sec_pos = self.parent.context_preserver._line_to_position(structural_analysis['raw_text'], section['line_number'])
9951:         distance = start_pos - sec_pos
9952:
9953:         # La sección debe preceder o estar en la unidad
9954:         if distance >= -50 and distance < min_distance:
9955:             best_match = section
9956:             min_distance = distance
9957:
9958:     if best_match:
9959:         hierarchy.append(best_match['title'])
9960:
9961:     return hierarchy
9962:
9963: def _enrich_strategic_metadata(self, units: List[Dict]) -> List[Dict]:
9964:     """Enriquecer unidades con metadatos estratégicos"""
9965:     for unit in units:
9966:         unit['strategic_weight'] = self._calculate_strategic_weight(unit)
9967:         unit['implementation_readiness'] = self._assess_implementation_readiness(unit)
9968:         unit['risk_level'] = self._assess_risk_level(unit)
9969:     return units
9970:
9971: def _calculate_strategic_weight(self, unit: Dict) -> float:
9972:     """Calcular peso estratégico de la unidad"""
9973:     factors = {
9974:         'chain_count': min(len(unit['chains']) / 5, 1.0),
9975:         'confidence': unit['confidence'],
9976:         'coherence': unit['semantic_coherence'],
9977:         'hierarchy_level': 1.0 if unit['hierarchy'] else 0.5
9978:     }
9979:     return np.mean(list(factors.values()))
9980:
9981: def _assess_implementation_readiness(self, unit: Dict) -> float:
9982:     """Evaluar preparación para implementación"""
9983:     text = unit.get('text', '')
9984:     readiness_indicators = [
9985:         'plan de acción', 'presupuesto asignado', 'cronograma definido',
9986:         'responsable designado', 'indicadores de seguimiento'
9987:     ]
9988:     readiness_score = sum(1 for ind in readiness_indicators if ind in text.lower())
9989:     return min(readiness_score / 3.0, 1.0)
9990:
9991: def _assess_risk_level(self, unit: Dict) -> float:
9992:     """Evaluar nivel de riesgo (simplificado)"""
9993:     text = unit.get('text', '')
9994:     risk_terms = ['riesgo', 'limitación', 'desafío', 'obstáculo', 'incertidumbre']
9995:     risk_score = sum(1 for term in risk_terms if term in text.lower())
9996:     return min(risk_score / 3.0, 1.0) #
9997:
9998:
9999: # =====
10000: # POLICY AREA CHUNK CALIBRATION - Garantiza 10 chunks por policy area
10001: # =====
10002:

```

```
10003: class PolicyAreaChunkCalibrator:
10004:     """
10005:         Calibrates chunking to guarantee exactly 10 strategic chunks per policy area.
10006:
10007:         Uses the existing SemanticChunkingProducer with dynamically adjusted parameters
10008:         to ensure consistent chunk count across different policy documents.
10009:
10010:     Strategy:
10011:         1. Estimate optimal chunk_size based on document length
10012:         2. Generate initial chunks with SemanticChunkingProducer
10013:         3. Adjust parameters iteratively if chunk count != 10
10014:         4. Merge or split chunks as needed to reach target
10015:
10016:     Attributes:
10017:         TARGET_CHUNKS_PER_PA: Target number of chunks (10)
10018:         TOLERANCE: Acceptable deviation (±1 chunk)
10019:         MAX_ITERATIONS: Maximum calibration iterations (3)
10020:     """
10021:
10022:     TARGET_CHUNKS_PER_PA = 10
10023:     TOLERANCE = 1
10024:     MAX_ITERATIONS = 3
10025:
10026:     # Canonical policy areas from questionnaire_monolith.json
10027:     POLICY_AREAS = [
10028:         "PA01", "PA02", "PA03", "PA04", "PA05",
10029:         "PA06", "PA07", "PA08", "PA09", "PA10"
10030:     ]
10031:
10032:     def __init__(self, semantic_chunking_producer: SemanticChunkingProducer):
10033:         """
10034:             Initialize calibrator.
10035:
10036:             Args:
10037:                 semantic_chunking_producer: Canonical SemanticChunkingProducer instance
10038:             """
10039:         self.chunking_producer = semantic_chunking_producer
10040:         self.logger = logging.getLogger("SPC.Calibrator")
10041:
10042:     def calibrate_for_policy_area(
10043:         self,
10044:         text: str,
10045:         policy_area: str,
10046:         metadata: Optional[Dict[str, Any]] = None
10047:     ) -> List[Dict[str, Any]]:
10048:         """
10049:             Generate exactly 10 chunks for a policy area.
10050:
10051:             Args:
10052:                 text: Policy document text
10053:                 policy_area: Policy area ID (PA01-PA10)
10054:                 metadata: Optional metadata to attach to chunks
10055:
10056:             Returns:
10057:                 List of exactly 10 chunks
10058:
```

```
10059:     Raises:
10060:         ValueError: If policy_area is invalid or text is empty
10061:         """
10062:     if policy_area not in self.POLICY_AREAS:
10063:         raise ValueError(
10064:             f"Invalid policy_area: {policy_area}. "
10065:             f"Must be one of {self.POLICY_AREAS}"
10066:         )
10067:
10068:     if not text or len(text.strip()) == 0:
10069:         raise ValueError("Text cannot be empty")
10070:
10071:     self.logger.info(
10072:         f"Calibrating chunks for {policy_area} "
10073:         f"(target: {self.TARGET CHUNKS_PER_PA} chunks)"
10074:     )
10075:
10076:     # Estimate initial parameters based on document length
10077:     initial_params = self._estimate_initial_params(text)
10078:
10079:     # Attempt to generate chunks with calibration
10080:     chunks = self._generate_with_calibration(
10081:         text,
10082:         policy_area,
10083:         initial_params,
10084:         metadata
10085:     )
10086:
10087:     # Final validation
10088:     if len(chunks) != self.TARGET_CHUNKS_PER_PA:
10089:         self.logger.warning(
10090:             f"{policy_area}: Could not reach exact target. "
10091:             f"Got {len(chunks)} chunks, forcing adjustment to {self.TARGET_CHUNKS_PER_PA}"
10092:         )
10093:         chunks = self._force_chunk_count(chunks, self.TARGET_CHUNKS_PER_PA)
10094:
10095:     self.logger.info(
10096:         f"{policy_area}: Calibration complete - {len(chunks)} chunks generated"
10097:     )
10098:
10099:     return chunks
10100:
10101: def _estimate_initial_params(self, text: str) -> Dict[str, Any]:
10102:     """
10103:     Estimate optimal chunking parameters based on text length.
10104:
10105:     Args:
10106:         text: Document text
10107:
10108:     Returns:
10109:         Dictionary with chunk_size, overlap, and other parameters
10110:         """
10111:     text_length = len(text)
10112:     sentence_count = text.count('.') + text.count('!') + text.count('?')
10113:
10114:     # Estimate chunk size to yield ~10 chunks
```

```

10115:         estimated_chunk_size = max(
10116:             500, # Minimum chunk size
10117:             min(
10118:                 2000, # Maximum chunk size
10119:                 text_length // (self.TARGET_CHUNKS_PER_PA + 2) # Add buffer
10120:             )
10121:         )
10122:
10123:     return {
10124:         'chunk_size': estimated_chunk_size,
10125:         'overlap': int(estimated_chunk_size * 0.15), # 15% overlap
10126:         'min_chunk_size': 300,
10127:         'adaptive': True,
10128:     }
10129:
10130:     def _generate_with_calibration(
10131:         self,
10132:         text: str,
10133:         policy_area: str,
10134:         initial_params: Dict[str, Any],
10135:         metadata: Optional[Dict[str, Any]]
10136:     ) -> List[Dict[str, Any]]:
10137:     """
10138:         Generate chunks with iterative calibration to reach target count.
10139:
10140:     Args:
10141:         text: Document text
10142:         policy_area: Policy area ID
10143:         initial_params: Initial chunking parameters
10144:         metadata: Optional metadata
10145:
10146:     Returns:
10147:         List of chunks (may not be exactly 10, needs final adjustment)
10148:     """
10149:     params = initial_params.copy()
10150:
10151:     for iteration in range(self.MAX_ITERATIONS):
10152:         # Use SemanticChunkingProducer to generate chunks
10153:         try:
10154:             # FIXED: Actually use the SemanticChunkingProducer
10155:             chunks = self._generate_chunks_with_producer(text, params, policy_area, metadata)
10156:         except Exception as e:
10157:             # Fallback to simple chunking if producer fails
10158:             self.logger.warning(f"SemanticChunkingProducer failed: {e}, using fallback")
10159:             chunks = self._generate_chunks_simple(text, params, policy_area, metadata)
10160:
10161:         chunk_count = len(chunks)
10162:         delta = chunk_count - self.TARGET_CHUNKS_PER_PA
10163:
10164:         self.logger.debug(
10165:             f"{policy_area}: Iteration {iteration+1} - "
10166:             f"{chunk_count} chunks (delta: {delta:+d})"
10167:         )
10168:
10169:         # Check if within tolerance
10170:         if abs(delta) <= self.TOLERANCE:

```

```

10171:             return chunks
10172:
10173:         # Adjust parameters for next iteration
10174:         if delta > 0:
10175:             # Too many chunks - increase chunk size
10176:             params['chunk_size'] = int(params['chunk_size'] * 1.2)
10177:         else:
10178:             # Too few chunks - decrease chunk size
10179:             params['chunk_size'] = int(params['chunk_size'] * 0.8)
10180:
10181:         # Ensure bounds
10182:         params['chunk_size'] = max(400, min(2500, params['chunk_size']))
10183:
10184:     # Return best attempt after max iterations
10185:     return chunks
10186:
10187:     def _generate_chunks_with_producer(
10188:         self,
10189:         text: str,
10190:         params: Dict[str, Any],
10191:         policy_area: str,
10192:         metadata: Optional[Dict[str, Any]]
10193:     ) -> List[Dict[str, Any]]:
10194:         """
10195:             Generate chunks using the SemanticChunkingProducer.
10196:
10197:             Args:
10198:                 text: Document text
10199:                 params: Chunking parameters (chunk_size, overlap, etc.)
10200:                 policy_area: Policy area ID
10201:                 metadata: Optional metadata
10202:
10203:             Returns:
10204:                 List of chunk dictionaries
10205:         """
10206:         # Use the actual SemanticChunkingProducer (instance method, not standalone function)
10207:         # Use the producer instance injected via __init__
10208:         producer = self.chunking_producer
10209:
10210:         # chunk_document signature: (text: str, preserve_structure: bool = True) -> list[dict[str, Any]]
10211:         result_chunks = producer.chunk_document(text=text, preserve_structure=True)
10212:
10213:         # Convert to our format
10214:         chunks = []
10215:         for i, chunk_result in enumerate(result_chunks):
10216:             # chunk_result is a dict with keys like 'text', 'embedding', 'section_type', etc.
10217:             chunks.append({
10218:                 'id': f"{policy_area}_chunk_{i+1}",
10219:                 'text': chunk_result.get('text', ''),
10220:                 'policy_area': policy_area,
10221:                 'chunk_index': i,
10222:                 'length': len(chunk_result.get('text', '')),
10223:                 'metadata': metadata or {},
10224:                 'semantic_metadata': {
10225:                     'section_type': chunk_result.get('section_type'),
10226:                     'section_id': chunk_result.get('section_id'),

```

```

10227:                 'has_embedding': 'embedding' in chunk_result
10228:             }
10229:         })
10230:
10231:     return chunks
10232:
10233: def _generate_chunks_simple(
10234:     self,
10235:     text: str,
10236:     params: Dict[str, Any],
10237:     policy_area: str,
10238:     metadata: Optional[Dict[str, Any]]
10239: ) -> List[Dict[str, Any]]:
10240:     """
10241:         Simple chunk generation using sentence splitting.
10242:
10243:     This is a fallback implementation. In production, this would use
10244:     the full SemanticChunkingProducer with BGE-M3 embeddings.
10245:
10246:     Args:
10247:         text: Document text
10248:         params: Chunking parameters
10249:         policy_area: Policy area ID
10250:         metadata: Optional metadata
10251:
10252:     Returns:
10253:         List of chunk dictionaries
10254:     """
10255:     # Simple sentence-based chunking
10256:     sentences = re.split(r'[.!?]+', text)
10257:     sentences = [s.strip() for s in sentences if s.strip()]
10258:
10259:     chunk_size = params.get('chunk_size', 1000)
10260:     chunks = []
10261:     current_chunk = []
10262:     current_length = 0
10263:
10264:     for sentence in sentences:
10265:         sentence_length = len(sentence)
10266:
10267:         if current_length + sentence_length > chunk_size and current_chunk:
10268:             # Create chunk
10269:             chunk_text = '. '.join(current_chunk) + '.'
10270:             chunks.append({
10271:                 'id': f'{policy_area}_chunk_{len(chunks)+1}',
10272:                 'text': chunk_text,
10273:                 'policy_area': policy_area,
10274:                 'chunk_index': len(chunks),
10275:                 'length': len(chunk_text),
10276:                 'metadata': metadata or {}
10277:             })
10278:             current_chunk = [sentence]
10279:             current_length = sentence_length
10280:         else:
10281:             current_chunk.append(sentence)
10282:             current_length += sentence_length

```

```

10283:
10284:     # Add final chunk
10285:     if current_chunk:
10286:         chunk_text = '. '.join(current_chunk) + '.'
10287:         chunks.append({
10288:             'id': f'{policy_area}_chunk_{len(chunks)+1}',
10289:             'text': chunk_text,
10290:             'policy_area': policy_area,
10291:             'chunk_index': len(chunks),
10292:             'length': len(chunk_text),
10293:             'metadata': metadata or {}
10294:         })
10295:
10296:     return chunks
10297:
10298: def _force_chunk_count(
10299:     self,
10300:     chunks: List[Dict[str, Any]],
10301:     target: int
10302: ) -> List[Dict[str, Any]]:
10303: """
10304:     Force chunk count to exactly match target by merging or splitting.
10305:
10306:     Args:
10307:         chunks: List of chunks
10308:         target: Target chunk count
10309:
10310:     Returns:
10311:         List with exactly target chunks
10312: """
10313: current_count = len(chunks)
10314:
10315: if current_count == target:
10316:     return chunks
10317:
10318: if current_count > target:
10319:     # Too many chunks - merge smallest adjacent pairs
10320:     while len(chunks) > target:
10321:         # Find smallest chunk
10322:         min_idx = min(range(len(chunks)), key=lambda i: chunks[i]['length'])
10323:
10324:         # Merge with adjacent chunk
10325:         if min_idx > 0:
10326:             # Merge with previous
10327:             chunks[min_idx-1]['text'] += ' ' + chunks[min_idx]['text']
10328:             chunks[min_idx-1]['length'] = len(chunks[min_idx-1]['text'])
10329:             chunks.pop(min_idx)
10330:         else:
10331:             # Merge with next
10332:             chunks[min_idx]['text'] += ' ' + chunks[min_idx+1]['text']
10333:             chunks[min_idx]['length'] = len(chunks[min_idx]['text'])
10334:             chunks.pop(min_idx+1)
10335:     else:
10336:         # Too few chunks - split largest chunks
10337:         while len(chunks) < target:
10338:             # Find largest chunk

```

```
10339:         max_idx = max(range(len(chunks)), key=lambda i: chunks[i]['length'])
10340:
10341:         # Split it in half
10342:         chunk_to_split = chunks[max_idx]
10343:         text = chunk_to_split['text']
10344:         mid_point = len(text) // 2
10345:
10346:         # Find sentence boundary near midpoint
10347:         split_point = text.rfind('.', 0, mid_point) + 1
10348:         if split_point <= 0:
10349:             split_point = mid_point
10350:
10351:         # Create two chunks
10352:         chunk1_text = text[:split_point].strip()
10353:         chunk2_text = text[split_point:].strip()
10354:
10355:         chunks[max_idx] = {
10356:             **chunk_to_split,
10357:             'text': chunk1_text,
10358:             'length': len(chunk1_text),
10359:         }
10360:
10361:         chunks.insert(max_idx + 1, {
10362:             **chunk_to_split,
10363:             'id': f"{chunk_to_split['id']}_split",
10364:             'text': chunk2_text,
10365:             'length': len(chunk2_text),
10366:         })
10367:
10368:     # Re-index chunks
10369:     for i, chunk in enumerate(chunks):
10370:         chunk['chunk_index'] = i
10371:
10372:     return chunks
10373:
10374:
10375: # =====
10376: # SISTEMA PRINCIPAL DE CHUNKING ESTRATÃ\211GICO (COMPLETO)
10377: # =====
10378:
10379: class StrategicChunkingSystem:
10380:     def __init__(self, random_seed: int = 42):
10381:         """
10382:             Initialize the Strategic Chunking System with canonical components.
10383:
10384:             Integrates production-grade canonical modules:
10385:             - PolicyAnalysisEmbedder for semantic embeddings
10386:             - SemanticProcessor for chunking with PDM structure awareness
10387:             - IndustrialPolicyProcessor for causal evidence extraction
10388:             - BayesianEvidenceScorer for probabilistic confidence scoring
10389:
10390:             Args:
10391:                 random_seed: Seed for deterministic RNG (default: 42)
10392:
10393:             Inputs:
10394:                 None
```

```

10395:     Outputs:
10396:         None - initializes system state
10397:         """
10398:         # Fix seeds for deterministic execution (HOSTILE AUDIT REQUIREMENT)
10399:         import random
10400:         np.random.seed(random_seed)
10401:         random.seed(random_seed)
10402:         self.logger = logging.getLogger("SPC") # Unified logger name
10403:         self.logger.info(f"Initialized with deterministic seed: {random_seed}")
10404:
10405:         self.config = SmartChunkConfig()
10406:
10407:         # =====
10408:         # CANONICAL SOTA PRODUCERS - Frontier approach components
10409:         # =====
10410:
10411:         # Initialize SOTA canonical producers that replace internal implementations
10412:         # These provide BGE-M3 embeddings, cross-encoder reranking, Bayesian numerical eval
10413:         self._spc_embed = EmbeddingPolicyProducer()           # chunking + embeddings + search + Bayesian numeric eval
10414:         self._spc_sem = SemanticChunkingProducer()          # direct embed_text/embed_batch and chunk_document
10415:         self._spc_policy = create_policy_processor()          # canonical PDQ/dimension evidence
10416:
10417:         self.logger.info("SOTA canonical producers initialized: EmbeddingPolicyProducer, SemanticChunkingProducer, PolicyProcessor")
10418:
10419:         # =====
10420:         # POLICY AREA \227 DIMENSION KEYWORD MAPS - Structured extraction
10421:         # =====
10422:
10423:         # PA01-PA10 keyword maps for content extraction
10424:         self._pa_keywords = {
10425:             "PA01": ["mujeres", "g@nero", "igualdad", "feminismo", "violencia g@nero", "empoderamiento", "equidad"],
10426:             "PA02": ["violencia", "conflicto armado", "protecci@n", "prevenci@n", "grupos delincuenciales", "econom@-as ilegales", "seguridad"],
10427:             "PA03": ["ambiente", "cambio clim@tico", "desastres", "medio ambiente", "ecolog@-a", "sostenibilidad", "recursos naturales"],
10428:             "PA04": ["derechos econ@micos", "derechos sociales", "derechos culturales", "educaci@n", "salud", "vivienda", "trabajo"],
10429:             "PA05": ["v@ctimas", "construcci@n de paz", "reconciliaci@n", "reparaci@n", "memoria", "justicia transicional"],
10430:             "PA06": ["ni@ez", "adolescencia", "juventud", "entornos protectores", "desarrollo infantil", "educaci@n inicial"],
10431:             "PA07": ["tierras", "territorios", "tenencia", "reforma agraria", "ordenamiento territorial", "catastro"],
10432:             "PA08": ["l@-deres", "lideres", "defensores", "defensoras", "derechos humanos", "protecci@n l@-deres", "amenazas"],
10433:             "PA09": ["privadas libertad", "c@rceles", "sistema penitenciario", "hacinamiento", "reinserci@n", "reclusos"],
10434:             "PA10": ["migraci@n", "transfronteriza", "migrantes", "refugiados", "movilidad humana", "frontera"]
10435:         }
10436:
10437:         # DIM01-DIM06 keyword maps for dimension alignment
10438:         self._dim_keywords = {
10439:             "DIM01": ["diagn@stico", "recursos", "presupuesto", "financiaci@n", "insumos", "inversi@n", "dotaci@n"],
10440:             "DIM02": ["actividades", "intervenci@n", "dise@to", "estrategias", "acciones", "programas", "proyectos"],
10441:             "DIM03": ["productos", "outputs", "entregables", "resultados intermedios", "metas", "indicadores producto"],
10442:             "DIM04": ["resultados", "outcomes", "efectos", "logros", "cambios", "impacto directo", "beneficiarios"],
10443:             "DIM05": ["impactos", "largo plazo", "transformaci@n", "cambio estructural", "sostenibilidad", "legado"],
10444:             "DIM06": ["causalidad", "teor@-a de cambio", "cadena causal", "l@-gica intervenci@n", "marco l@-gico", "supuestos"]
10445:         }
10446:
10447:         self.logger.info(f"PA keyword maps: {len(self._pa_keywords)} policy areas")
10448:         self.logger.info(f"DIM keyword maps: {len(self._dim_keywords)} dimensions")
10449:
10450:         # =====

```

```

10451:     # SPECIALIZED COMPONENTS - Keep (no canonical equivalent)
10452:     # =====
10453:
10454:     # These provide unique Smart Policy Chunks innovations
10455:     self._nlp = None # SpaCy for NER (lazy-loaded)
10456:     self._kg_builder = None # NetworkX knowledge graph
10457:     self._topic_modeler = None # LDA topic modeling
10458:     self._argument_analyzer = None # Toulmin argument structure
10459:     self._temporal_analyzer = None # Temporal dynamics
10460:     self._discourse_analyzer = None # Discourse markers
10461:     self._strategic_integrator = None # Cross-reference integration
10462:     self._causal_analyzer = None # Causal chain analyzer (lazy-loaded)
10463:
10464:     # Modelo para clasificaciÃ³n de tipo de chunk
10465:     self.chunk_classifier = None
10466:
10467:     # Almacenamiento
10468:     self.tfidf_vectorizer = TfidfVectorizer(stop_words=self._get_stopwords(), ngram_range=(1, 2), max_df=0.85, min_df=2)
10469:     self.chunks_for_tfidf = []
10470:     self.corpus_embeddings = None
10471:
10472:     # =====
10473:     # SPECIALIZED COMPONENT PROPERTIES - Innovation layers (no canonical equivalent)
10474:     #
10475:
10476:     @property
10477:     def nlp(self):
10478:         """Lazy-load SpaCy NLP model"""
10479:         if self._nlp is None:
10480:             try:
10481:                 self.logger.info("Loading SpaCy model: es_core_news_lg")
10482:                 self._nlp = spacy.load("es_core_news_lg")
10483:             except:
10484:                 self.logger.warning("SpaCy es_core_news_lg no disponible, usando sm")
10485:             try:
10486:                 self._nlp = spacy.load("es_core_news_sm")
10487:             except:
10488:                 self.logger.error("NingÃ³n modelo SpaCy disponible. Funcionalidad de NER limitada.")
10489:                 self._nlp = None
10490:             return self._nlp
10491:
10492:     @property
10493:     def context_preserver(self):
10494:         """
10495:             CANONICAL REPLACEMENT: Use semantic_processor for chunking.
10496:             Kept for backward compatibility but delegates to canonical component.
10497:         """
10498:         # Return a lightweight adapter that uses canonical SemanticProcessor
10499:         return self.semantic_processor
10500:
10501:     @property
10502:     def causal_analyzer(self):
10503:         """Lazy-load causal chain analyzer"""
10504:         if self._causal_analyzer is None:
10505:             self._causal_analyzer = CausalChainAnalyzer(self)
10506:         return self._causal_analyzer

```

```
10507:  
10508:     @property  
10509:     def kg_builder(self):  
10510:         """Lazy-load knowledge graph builder"""  
10511:         if self._kg_builder is None:  
10512:             self._kg_builder = KnowledgeGraphBuilder(self)  
10513:         return self._kg_builder  
10514:  
10515:     @property  
10516:     def topic_modeler(self):  
10517:         """Lazy-load topic modeler"""  
10518:         if self._topic_modeler is None:  
10519:             self._topic_modeler = TopicModeler(self)  
10520:         return self._topic_modeler  
10521:  
10522:     @property  
10523:     def argument_analyzer(self):  
10524:         """Lazy-load argument analyzer"""  
10525:         if self._argument_analyzer is None:  
10526:             self._argument_analyzer = ArgumentAnalyzer(self)  
10527:         return self._argument_analyzer  
10528:  
10529:     @property  
10530:     def temporal_analyzer(self):  
10531:         """Lazy-load temporal analyzer"""  
10532:         if self._temporal_analyzer is None:  
10533:             self._temporal_analyzer = TemporalAnalyzer(self)  
10534:         return self._temporal_analyzer  
10535:  
10536:     @property  
10537:     def discourse_analyzer(self):  
10538:         """Lazy-load discourse analyzer"""  
10539:         if self._discourse_analyzer is None:  
10540:             self._discourse_analyzer = DiscourseAnalyzer(self)  
10541:         return self._discourse_analyzer  
10542:  
10543:     @property  
10544:     def strategic_integrator(self):  
10545:         """Lazy-load strategic integrator"""  
10546:         if self._strategic_integrator is None:  
10547:             self._strategic_integrator = StrategicIntegrator(self)  
10548:         return self._strategic_integrator  
10549:  
10550:     def detect_language(self, text: str) -> str:  
10551:         """  
10552:             Detect the primary language of a text document.  
10553:  
10554:             Inputs:  
10555:                 text (str): Text to analyze  
10556:             Outputs:  
10557:                 str: ISO 639-1 language code (e.g., 'es', 'en', 'pt')  
10558:         """  
10559:         if not LANGDETECT_AVAILABLE:  
10560:             # Default to Spanish for Colombian policy documents  
10561:             return 'es'  
10562:
```

```

10563:     try:
10564:         # Sample first 2000 characters for language detection
10565:         sample = safe_utf8_truncate(text, 2000)
10566:         detected_lang = detect(sample)
10567:         self.logger.info(f"Detected language: {detected_lang}")
10568:         return detected_lang
10569:     except Exception as e:
10570:         self.logger.warning(f"Language detection failed: {e}, defaulting to Spanish")
10571:         return 'es'
10572:
10573:     def select_embedding_model_for_language(self, language: str) -> None:
10574:         """
10575:             Select appropriate embedding model based on detected language.
10576:
10577:             Inputs:
10578:                 language (str): ISO 639-1 language code
10579:             Outputs:
10580:                 None - updates model selection
10581:         """
10582:         # For now, multilingual-e5-large handles multiple languages well
10583:         # Could be extended with language-specific models if needed
10584:         if language in ['es', 'pt', 'ca']: # Spanish, Portuguese, Catalan
10585:             self.logger.info(f"Using multilingual model for {language} (optimal for Romance languages)")
10586:         else:
10587:             self.logger.info(f"Using multilingual model for {language}")
10588:
10589:         # Model is already multilingual, no change needed
10590:         # This method provides extension point for future language-specific optimization
10591:
10592:     # --- MÃ©todos de la clase principal (ContinuaciÃ³n de smart_policy_chunks_industrial_v3_complete_Version2.py) ---
10593:
10594:     def _get_stopwords(self) -> List[str]:
10595:         """
10596:             Get Spanish stopwords list.
10597:
10598:             Inputs:
10599:                 None
10600:             Outputs:
10601:                 List[str]: List of Spanish stopwords
10602:         """
10603:         # Una lista de stopwords mÃ¡s completa se usarÃ¡ en producciÃ³n
10604:         return ['el', 'la', 'los', 'las', 'un', 'una', 'unos', 'unas', 'y', 'o', 'de', 'a', 'en', 'por', 'con', 'para', 'del', 'al', 'que', 'se', 'es', 'son',
10605: 'han', 'como', 'mÃ¡s', 'pero', 'no', 'su', 'sus', 'ha', 'lo', 'e', 'u', 'ni', 'sin', 'mi', 'tu', 'si', 'cuando', 'este', 'esta', 'estos', 'estas', 'esos', 'esas',
10606: 'aquel', 'aquella']
10607:     # =====
10608:     # CANONICAL SOTA METHODS - Replace manual implementations
10609:     # =====
10610:     def semantic_search_with_rerank(
10611:         self,
10612:         query: str,
10613:         chunks: list[dict],
10614:         pdq_filter: dict | None = None,
10615:         top_k: int = 10
10616:     ) -> list[tuple[dict, float]]:

```

```
10617:      """
10618:      CANONICAL SOTA: Semantic search with cross-encoder reranking.
10619:
10620:      Replaces manual cosine_similarity ranking with SOTA reranker.
10621:
10622:      Inputs:
10623:          query (str): Search query
10624:          chunks (list[dict]): Chunks to search
10625:          pdq_filter (dict | None): Optional PDQ filter
10626:          top_k (int): Number of results to return
10627:      Outputs:
10628:          list[tuple[dict, float]]: List of (chunk, score) tuples
10629:      """
10630:      results = self._spc_embed.semantic_search(
10631:          query, chunks, pdq_filter=pdq_filter, use_reranking=True
10632:      )
10633:      return results[:top_k]
10634:
10635:  def _attach_canonical_evidence(self, full_text: str) -> dict[str, Any]:
10636:      """
10637:          CANONICAL SOTA: Attach canonical PDQ/dimension evidence.
10638:
10639:          Uses canonical policy patterns instead of ad-hoc heuristics.
10640:
10641:          Inputs:
10642:              full_text (str): Full document text
10643:          Outputs:
10644:              dict[str, Any]: Canonical evidence analysis
10645:      """
10646:      return self._spc_policy.analyze_text(full_text)
10647:
10648:  def evaluate_numerical_consistency(
10649:      self,
10650:      chunks: list[dict],
10651:      pdq_context: dict
10652:  ) -> dict[str, Any]:
10653:      """
10654:          CANONICAL SOTA: Evaluate numerical consistency with Bayesian analysis.
10655:
10656:          Uses canonical extractor + Bayesian analyzer for probabilistic scoring.
10657:
10658:          Inputs:
10659:              chunks (list[dict]): Chunks to evaluate
10660:              pdq_context (dict): PDQ context for evaluation
10661:          Outputs:
10662:              dict[str, Any]: Numerical consistency evaluation
10663:      """
10664:      return self._spc_embed.evaluate_numerical_consistency(chunks, pdq_context)
10665:
10666:  def _generate_embedding(self, text: str, model_type: str = "semantic") -> np.ndarray:
10667:      """
10668:          CANONICAL SOTA: Generate embedding using SemanticChunkingProducer.
10669:
10670:          Replaces internal embedding with SOTA multilingual BGE-M3 model.
10671:          model_type kept for compatibility; canonical pipeline is multilingual.
10672:
```

```
10673:     Inputs:
10674:         text (str): Input text to embed
10675:         model_type (str): Ignored, canonical uses SOTA multilingual
10676:     Outputs:
10677:         np.ndarray: Embedding vector from canonical SOTA component
10678:     """
10679:     return self._spc_sem.embed_text(text).astype(np.float32)
10680:
10681:     def _create_smart_policy_chunk(
10682:         self,
10683:         strategic_unit: Dict,
10684:         metadata: Dict,
10685:         argument_structure: Optional[ArgumentStructure],
10686:         temporal_structure: Optional[TemporalDynamics],
10687:         discourse_structure: Optional[Dict],
10688:         global_topics: Dict,
10689:         global_kg: Dict
10690:     ) -> SmartPolicyChunk:
10691:         """Crear Smart Policy Chunk con análisis completo"""
10692:         text = strategic_unit.get("text", "")
10693:         document_id = metadata.get("document_id", "doc_001")
10694:
10695:         # Generar embeddings
10696:         semantic_embedding = self._generate_embedding(text, 'semantic')
10697:         policy_embedding = self._generate_embedding(text, 'semantic')
10698:         causal_embedding = self._generate_embedding(text, 'semantic')
10699:         temporal_embedding = self._generate_embedding(text, 'semantic')
10700:
10701:         # Análisis causales
10702:         causal_evidence = self._extract_comprehensive_causal_evidence(strategic_unit, global_kg)
10703:
10704:         # Entidades políticas
10705:         policy_entities = self._extract_policy_entities_with_context(strategic_unit)
10706:
10707:         # Contexto estratégico
10708:         strategic_context = self._derive_strategic_context(strategic_unit, global_topics)
10709:
10710:         # MÁTRICAS de calidad
10711:         confidence_metrics = self._calculate_comprehensive_confidence(strategic_unit, causal_evidence, policy_entities)
10712:         coherence_score = self._calculate_coherence_score(strategic_unit)
10713:         completeness_index = self._calculate_completeness_index(strategic_unit, causal_evidence)
10714:         strategic_importance = self._assess_strategic_importance(strategic_unit, causal_evidence, policy_entities)
10715:         information_density = self._calculate_information_density(text)
10716:         actionability_score = self._assess_actionability(text, policy_entities)
10717:
10718:         # Referencias cruzadas
10719:         cross_refs = self._find_cross_document_references(strategic_unit)
10720:
10721:         # Supuestos implícitos
10722:         implicit_assumptions = self._extract_implicit_assumptions(strategic_unit, causal_evidence)
10723:
10724:         # Presuposiciones contextuales
10725:         contextual_presuppositions = self._identify_contextual_presuppositions(strategic_unit)
10726:
10727:         # Marcadores de discurso
10728:         discourse_markers = self._extract_discourse_markers(text)
```

```
10729:  
10730:      # Patrones retóricos  
10731:      rhetorical_patterns = self._identify_rhetorical_patterns(text)  
10732:  
10733:      # Distribución de tópicos para este chunk  
10734:      topic_distribution = self._calculate_chunk_topic_distribution(text, global_topics)  
10735:  
10736:      # Frases clave  
10737:      key_phrases = self._extract_key_phrases(text)  
10738:  
10739:      # Nodos y aristas del grafo de conocimiento  
10740:      kg_nodes = [e.normalized_form for e in policy_entities]  
10741:      kg_edges = self._derive_kg_edges_for_chunk(policy_entities, causal_evidence)  
10742:  
10743:      # Hash y ID  
10744:      content_hash = hashlib.sha256(text.encode('utf-8')).hexdigest()  
10745:      chunk_id = f"{document_id}_{content_hash[:8]}"  
10746:  
10747:      # Normalización  
10748:      normalized_text = self._advanced_preprocessing(text)  
10749:  
10750:      return SmartPolicyChunk(  
10751:          chunk_id=chunk_id,  
10752:          document_id=document_id,  
10753:          content_hash=content_hash,  
10754:          policy_area_id=strategic_unit.get("policy_area_id"), # PA01-PA10  
10755:          dimension_id=strategic_unit.get("dimension_id"), # DIM01-DIM06  
10756:  
10757:          text=text,  
10758:          normalized_text=normalized_text,  
10759:          semantic_density=self._calculate_semantic_density(text),  
10760:  
10761:          section_hierarchy=strategic_unit.get("hierarchy", []),  
10762:          document_position=strategic_unit.get("position", (0, 0)),  
10763:          chunk_type=self._classify_chunk_type(text),  
10764:  
10765:          causal_chain=causal_evidence,  
10766:          policy_entities=policy_entities,  
10767:          implicit_assumptions=implicit_assumptions,  
10768:          contextual_presuppositions=contextual_presuppositions,  
10769:  
10770:          argument_structure=argument_structure,  
10771:          temporal_dynamics=temporal_structure,  
10772:          discourse_markers=discourse_markers,  
10773:          rhetorical_patterns=rhetorical_patterns,  
10774:  
10775:          cross_references=cross_refs,  
10776:          strategic_context=strategic_context,  
10777:  
10778:          confidence_metrics=confidence_metrics,  
10779:          coherence_score=coherence_score,  
10780:          completeness_index=completeness_index,  
10781:          strategic_importance=strategic_importance,  
10782:          information_density=information_density,  
10783:          actionability_score=actionability_score,  
10784:
```

```
10785:         semantic_embedding=semantic_embedding,
10786:         policy_embedding=policy_embedding,
10787:         causal_embedding=causal_embedding,
10788:         temporal_embedding=temporal_embedding,
10789:
10790:         knowledge_graph_nodes=kg_nodes,
10791:         knowledge_graph_edges=kg_edges,
10792:
10793:         topic_distribution=topic_distribution,
10794:         key_phrases=key_phrases,
10795:
10796:         model_versions=self._get_model_versions()
10797:     )
10798:
10799: # --- MÃ©todos de AnÃ;lisis y ExtracciÃ³n (ContinuaciÃ³n de smart_policy_chunks_industrial_v3_complete_Version2.py) ---
10800:
10801: def _extract_comprehensive_causal_evidence(self, strategic_unit: Dict, global_kg: Dict) -> List[CausalEvidence]:
10802:     """Extraer evidencia causal completa (Placeholder con estructura avanzada)"""
10803:     evidence_list = []
10804:     text = strategic_unit.get("text", "")
10805:
10806:     # Placeholder para un modelo de extracciÃ³n causal mÃ;s avanzado
10807:     # SimulaciÃ³n de extracciÃ³n basada en patrones de palabras
10808:
10809:     dimensions = {
10810:         'problem_solution': [r'soluciÃ³n\s+para\s+([^.]+)', r'abordar\s+([^.]+)\s+mediante\s+([^.]+)'],
10811:         'impact_assessment': [r'tendrÃ; impacto\s+([^.]+)', r'conducirÃ; \s+a\s+([^.]+)'],
10812:         'resource_allocation': [r'asignaciÃ³n\s+de\s+([^.]+)\s+para\s+([^.]+)'],
10813:         'policy_instrument': [r'(?:la\s+implementaciÃ³n|el\s+uso)\s+de\s+([^.]+)\s+resultarÃ; \s+en\s+([^.]+)']
10814:     }
10815:
10816:     for dimension, patterns in dimensions.items():
10817:         for pattern in patterns:
10818:             matches = list(re.finditer(pattern, text, re.IGNORECASE))
10819:             for match in matches:
10820:                 context_start = max(0, match.start() - 150)
10821:                 context_end = min(len(text), match.end() + 150)
10822:                 context = text[context_start:context_end]
10823:
10824:                 # Determinar tipo de relaciÃ³n causal
10825:                 causal_type = self._determine_causal_type(match.group(0), context)
10826:
10827:                 # Calcular fuerza de la relaciÃ³n
10828:                 strength = self._calculate_causal_strength(match.group(0), context)
10829:
10830:                 # Identificar mecanismos
10831:                 mechanisms = self._identify_causal_mechanisms(context)
10832:
10833:                 evidence_list.append(CausalEvidence(
10834:                     dimension=dimension,
10835:                     category=self._categorize_causal_evidence(dimension),
10836:                     matches=[match.group(0)],
10837:                     confidence=0.75 + (strength * 0.2), # Ajuste por fuerza
10838:                     context_span=(context_start, context_end),
10839:                     implicit_indicators=self._find_implicit_indicators(context),
10840:                     causal_type=causal_type,
```

```
10841:             strength_score=strength,
10842:             mechanisms=mechanisms,
10843:             confounders=self._identify_confounders(context),
10844:             mediators=self._identify_mediators(context),
10845:             moderators=self._identify_moderators(context)
10846:         ))
10847:
10848:     return evidence_list
10849:
10850: def _categorize_causal_evidence(self, dimension: str) -> str:
10851:     """Categorizar el tipo de evidencia causal"""
10852:     if dimension in ['problem_solution', 'impact_assessment']:
10853:         return 'macro_policy'
10854:     elif dimension in ['resource_allocation', 'policy_instrument']:
10855:         return 'implementation_mechanisms'
10856:     return 'general'
10857:
10858: def _determine_causal_type(self, match_text: str, context: str) -> CausalRelationType:
10859:     """Determinar el tipo de relaciÃ³n causal por marcadores lingÃ¼Ã¢sticos"""
10860:     match_lower = match_text.lower()
10861:     if 'si' in match_lower and 'entonces' in context.lower():
10862:         return CausalRelationType.CONDITIONAL
10863:     elif any(word in match_lower for word in ['porque', 'debido a', 'gracias a']):
10864:         return CausalRelationType.DIRECT_CAUSE
10865:     elif any(word in match_lower for word in ['por lo tanto', 'en consecuencia']):
10866:         return CausalRelationType.DIRECT_CAUSE
10867:     else:
10868:         return CausalRelationType.INDIRECT_CAUSE
10869:
10870: def _calculate_causal_strength(self, match_text: str, context: str) -> float:
10871:     """Calcular fuerza de relaciÃ³n causal"""
10872:     strength_score = 0.5
10873:
10874:     # Reforzadores (Boosters)
10875:     boosters = ['directamente', 'significativamente', 'claramente', 'contundentemente']
10876:     if any(b in context.lower() for b in boosters):
10877:         strength_score += 0.3
10878:
10879:     # Mitigadores (Dampeners)
10880:     dampeners = ['parcialmente', 'limitadamente', 'posiblemente', 'podrÃ¡-a']
10881:     if any(d in context.lower() for d in dampeners):
10882:         strength_score -= 0.3
10883:
10884:     return np.clip(strength_score, 0.1, 1.0)
10885:
10886: def _identify_causal_mechanisms(self, context: str) -> List[str]:
10887:     """Identificar mecanismos causales (cÃ³mo funciona la relaciÃ³n)"""
10888:     mechanisms = []
10889:     mechanism_patterns = [
10890:         r'a\s+travÃ©s\s+de\s+([^.]+)', 
10891:         r'mediante\s+([^.]+)'
10892:     ]
10893:
10894:     for pattern in mechanism_patterns:
10895:         matches = re.findall(pattern, context, re.IGNORECASE)
10896:         mechanisms.extend(matches[:2])
```

```
10897:  
10898:         return mechanisms  
10899:  
10900:     def _find_implicit_indicators(self, context: str) -> List[str]:  
10901:         """Encontrar indicadores implÃ¡citos de causalidad (sin marcadores explÃ¡citos)"""  
10902:         implicit = []  
10903:         implicit_patterns = [  
10904:             r'para\s+lograr\s+([^.]+)',  
10905:             r'es\s+fundamental\s+([^.]+)',  
10906:             r'la\s+falta\s+de\s+([^.]+)\s+impacta\s+([^.]+)'  
10907:         ]  
10908:  
10909:         for pattern in implicit_patterns:  
10910:             matches = re.findall(pattern, context, re.IGNORECASE)  
10911:             implicit.extend(matches[:2])  
10912:  
10913:         return implicit  
10914:  
10915:     def _identify_confounders(self, context: str) -> List[str]:  
10916:         """Identificar factores de confusiÃ³n causales"""  
10917:         confounders = []  
10918:         confounder_patterns = [  
10919:             r'a\s+pesar\s+de\s+([^.]+)',  
10920:             r'sin\s+considerar\s+([^.]+)'  
10921:         ]  
10922:  
10923:         for pattern in confounder_patterns:  
10924:             matches = re.findall(pattern, context, re.IGNORECASE)  
10925:             confounders.extend(matches[:2])  
10926:  
10927:         return confounders  
10928:  
10929:     def _identify_mediators(self, context: str) -> List[str]:  
10930:         """Identificar mediadores causales"""  
10931:         mediators = []  
10932:         mediator_patterns = [  
10933:             r'que\s+a\s+su\s+vez\s+([^.]+)',  
10934:             r'lo\s+cual\s+([^.]+)'  
10935:         ]  
10936:  
10937:         for pattern in mediator_patterns:  
10938:             matches = re.findall(pattern, context, re.IGNORECASE)  
10939:             mediators.extend(matches[:2])  
10940:  
10941:         return mediators  
10942:  
10943:     def _identify_moderators(self, context: str) -> List[str]:  
10944:         """Identificar moderadores causales"""  
10945:         moderators = []  
10946:         moderator_patterns = [  
10947:             r'en\s+la\s+medida\s+(?:en\s+)?que\s+([^.]+)',  
10948:             r'siempre\s+y\s+cuando\s+([^.]+)',  
10949:             r'dependiendo\s+de\s+([^.]+)'  
10950:         ]  
10951:  
10952:         for pattern in moderator_patterns:
```

```

10953:         matches = re.findall(pattern, context, re.IGNORECASE)
10954:         moderators.extend(matches[:2])
10955:
10956:     return moderators
10957:
10958: def _extract_policy_entities_with_context(self, strategic_unit: Dict) -> List[PolicyEntity]:
10959:     """Extraer entidades de polÃtica con contexto completo"""
10960:     entities = []
10961:     text = strategic_unit.get("text", "")
10962:
10963:     # Patrones de entidades institucionales
10964:     institutional_patterns = [
10965:         (r'(?:Ministerio|SecretarÃ-a|Departamento|DirecciÃ³n|Instituto|Agencia)\s+(?:de|del?|para)\s+[A-ZÃ±\201Ã±\211Ã±\215Ã±\223Ã±\232Ã±\221][a-zÃ±;Ã©Ã±-Ã³Ã°Ã±\s]+', 'institution'),
10966:         (r'(?:AlcaldÃ-a|GobernaciÃ³n|Prefectura)\s+(?:de|del?)\s+[A-ZÃ±\201Ã±\211Ã±\215Ã±\223Ã±\232Ã±\221][a-zÃ±;Ã©Ã±-Ã³Ã°Ã±\s]+', 'local_government'),
10967:         (r'(?:Consejo|ComitÃ©|Junta)\s+(?:de|del?|para)\s+[A-ZÃ±\201Ã±\211Ã±\215Ã±\223Ã±\232Ã±\221][a-zÃ±;Ã©Ã±-Ã³Ã°Ã±\s]+', 'committee')
10968:     ]
10969:
10970:     for pattern, entity_type in institutional_patterns:
10971:         matches = re.finditer(pattern, text)
10972:         for match in matches:
10973:             role = self._infer_entity_role(match.group(0), text)
10974:             entities.append(PolicyEntity(
10975:                 entity_type=entity_type,
10976:                 text=match.group(0),
10977:                 normalized_form=self._normalize_entity(match.group(0)),
10978:                 context_role=role,
10979:                 confidence=0.7,
10980:                 span=match.span()
10981:             ))
10982:
10983:     # Patrones de poblaciÃ³n/beneficiarios (simplificado)
10984:     beneficiary_patterns = [
10985:         (r'(?:los|las)\s+(?:niÃ±os|niÃ±as|jÃ³venes|mujeres|poblaciÃ³n|comunidad|ciudadanos)\s+(?:de|en)\s+[^,.]+', 'population_group'),
10986:         (r'beneficiarios\s+(?:de|del?)\s+[^,.]+', 'beneficiary_group')
10987:     ]
10988:
10989:     for pattern, entity_type in beneficiary_patterns:
10990:         matches = re.finditer(pattern, text, re.IGNORECASE)
10991:         for match in matches:
10992:             entities.append(PolicyEntity(
10993:                 entity_type=entity_type,
10994:                 text=match.group(0),
10995:                 normalized_form=self._normalize_entity(match.group(0)),
10996:                 context_role=PolicyEntityRole.BENEFICIARY,
10997:                 confidence=0.8,
10998:                 span=match.span()
10999:             ))
11000:
11001:     # DeduplicaciÃ³n y conteo
11002:     unique_entities: List[PolicyEntity] = []
11003:     seen_texts = set()
11004:     for entity in entities:
11005:         normalized = entity.normalized_form.lower()
11006:         if normalized not in seen_texts:
11007:             seen_texts.add(normalized)

```

```

11008:         unique_entities.append(entity)
11009:     else:
11010:         # Actualizar contador de menciones
11011:         for existing in unique_entities:
11012:             if existing.normalized_form.lower() == normalized:
11013:                 existing.mentioned_count += 1
11014:                 break
11015:
11016:     return unique_entities
11017:
11018: def _infer_entity_role(self, entity_text: str, context: str) -> PolicyEntityRole:
11019:     """Inferir el rol de la entidad en el contexto"""
11020:     context_lower = context.lower()
11021:     if any(term in context_lower for term in [f"ejecutarÃ; {entity_text.lower()}", f"responsable de {entity_text.lower()}", f"{entity_text.lower()} implementarÃ;"]):
11022:         return PolicyEntityRole.EXECUTOR
11023:     elif any(term in context_lower for term in [f"beneficiarÃ; a {entity_text.lower()}", f"dirigido a {entity_text.lower()}"]):
11024:         return PolicyEntityRole.BENEFICIARY
11025:     elif any(term in context_lower for term in [f"regulador {entity_text.lower()}", f"{entity_text.lower()} emitirÃ;"]):
11026:         return PolicyEntityRole.REGULATOR
11027:     elif any(term in context_lower for term in [f"financiarÃ; {entity_text.lower()}", f"funder {entity_text.lower()}"]):
11028:         return PolicyEntityRole.FUNDER
11029:     else:
11030:         return PolicyEntityRole.STAKEHOLDER
11031:
11032: def _normalize_entity(self, entity_text: str) -> str:
11033:     """Normalizar la forma de la entidad"""
11034:     # Eliminar artÃculos, preposiciones y estandarizar mayÃsculas
11035:     text = re.sub(r'\b(el|la|los|las|un|una|de|del|a|en|por)\b', '', entity_text, flags=re.IGNORECASE).strip()
11036:     text = re.sub(r'\s+', ' ', text)
11037:     return text.title()
11038:
11039: def _derive_strategic_context(self, strategic_unit: Dict, global_topics: Dict) -> StrategicContext:
11040:     """Derivar contexto estratÃgico comprensivo"""
11041:     text = strategic_unit.get("text", "")
11042:
11043:     return StrategicContext(
11044:         policy_intent=self._infer_policy_intent(text),
11045:         implementation_phase=self._identify_implementation_phase(text),
11046:         geographic_scope=self._extract_geographic_scope(text),
11047:         temporal_horizon=self._determine_temporal_horizon(text),
11048:         budget_linkage=self._identify_budget_linkage(text),
11049:         risk_factors=self._extract_risk_factors(text),
11050:         success_indicators=self._identify_success_indicators(text),
11051:         alignment_with_sdg=self._identify_sdg_alignment(text),
11052:         stakeholder_map=self._build_stakeholder_map(text),
11053:         policy_coherence_score=self._calculate_policy_coherence(text, global_topics),
11054:         intervention_logic_chain=self._extract_intervention_logic(strategic_unit)
11055:     )
11056:
11057: def _infer_policy_intent(self, text: str) -> str:
11058:     """Inferir la intenciÃn de polÃtica (e.g., mitigar, promover, regular)"""
11059:     text_lower = text.lower()
11060:     if 'promover' in text_lower or 'fomentar' in text_lower or 'impulsar' in text_lower:
11061:         return 'PromociÃn/Fomento'
11062:     elif 'reducir' in text_lower or 'mitigar' in text_lower or 'combatir' in text_lower:

```

```
11063:         return 'MitigaciÃ³n/ReducciÃ³n'
11064:     elif 'regular' in text_lower or 'establecer normativa' in text_lower or 'ley' in text_lower:
11065:         return 'RegulaciÃ³n/Normativa'
11066:     return 'General'
11067:
11068: def _identify_implementation_phase(self, text: str) -> str:
11069:     """Identificar la fase de implementaciÃ³n (e.g., diseÃ±o, ejecuciÃ³n, evaluaciÃ³n)"""
11070:     text_lower = text.lower()
11071:     if 'diseÃ±o' in text_lower or 'formulaciÃ³n' in text_lower:
11072:         return 'DiseÃ±o'
11073:     elif 'ejecuciÃ³n' in text_lower or 'implementaciÃ³n' in text_lower or 'puesta en marcha' in text_lower:
11074:         return 'EjecuciÃ³n'
11075:     elif 'evaluaciÃ³n' in text_lower or 'seguimiento' in text_lower or 'monitoreo' in text_lower:
11076:         return 'EvaluaciÃ³n'
11077:     return 'Mixto'
11078:
11079: def _extract_geographic_scope(self, text: str) -> str:
11080:     """Extraer el alcance geogrÃ;fico"""
11081:     text_lower = text.lower()
11082:     if 'nacional' in text_lower or 'paÃ;s' in text_lower:
11083:         return 'Nacional'
11084:     elif 'departamental' in text_lower or 'departamento' in text_lower or 'provincia' in text_lower:
11085:         return 'Departamental'
11086:     elif 'municipal' in text_lower or 'municipio' in text_lower or 'ciudad' in text_lower:
11087:         return 'Municipal/Local'
11088:     return 'Indefinido'
11089:
11090: def _determine_temporal_horizon(self, text: str) -> str:
11091:     """Determinar el horizonte temporal (corto, mediano, largo plazo)"""
11092:     text_lower = text.lower()
11093:     if 'corto plazo' in text_lower or 'prÃ³ximo aÃ±o' in text_lower:
11094:         return 'Corto Plazo'
11095:     elif 'mediano plazo' in text_lower or 'prÃ³ximos 4 aÃ±os' in text_lower:
11096:         return 'Mediano Plazo'
11097:     elif 'largo plazo' in text_lower or 'horizonte 2030' in text_lower:
11098:         return 'Largo Plazo'
11099:     return 'Mixto'
11100:
11101: def _identify_budget_linkage(self, text: str) -> str:
11102:     """Identificar vÃ;nculos presupuestales"""
11103:     text_lower = text.lower()
11104:     if 'presupuesto' in text_lower or 'recursos financieros' in text_lower or 'inversiÃ³n de' in text_lower:
11105:         return 'ExplÃ-cito'
11106:     elif 'costos' in text_lower or 'financiamiento' in text_lower:
11107:         return 'ImplÃ-cito'
11108:     return 'Ausente'
11109:
11110: def _extract_risk_factors(self, text: str) -> List[str]:
11111:     """Extraer factores de riesgo explÃ-citos"""
11112:     risks = []
11113:     risk_patterns = [
11114:         r'el\s+riesgo\s+(?:de|es)\s+([^.]+)', 
11115:         r'se\s+deben\s+mitigar\s+([^.]+)'
11116:     ]
11117:     for pattern in risk_patterns:
11118:         matches = re.findall(pattern, text, re.IGNORECASE)
```

```

11119:         risks.extend(matches[:3])
11120:     return risks
11121:
11122:     def _identify_success_indicators(self, text: str) -> List[str]:
11123:         """Identificar indicadores de éxito o máticas"""
11124:         indicators = []
11125:         indicator_patterns = [
11126:             r'indicador\s+([^.]+)',
11127:             r'meta\s+(?:de|para)\s+([^.]+)',
11128:             r'se\s+medirÁ;s+con\s+([^.]+)'
11129:         ]
11130:         for pattern in indicator_patterns:
11131:             matches = re.findall(pattern, text, re.IGNORECASE)
11132:             indicators.extend(matches[:3])
11133:     return indicators
11134:
11135:     def _identify_sdg_alignment(self, text: str) -> List[str]:
11136:         """Identificar alineación con ODS"""
11137:         sdgs = []
11138:         # ODS explícitos
11139:         sdg_pattern = r'ODS\s+(\d+\''
11140:         matches = re.findall(sdg_pattern, text, re.IGNORECASE)
11141:         sdgs.extend([f"ODS_{m}" for m in matches])
11142:
11143:         # ODS por temas (simplificado)
11144:         sdg_themes = {
11145:             'ODS_1': ['pobreza', 'pobres'],
11146:             'ODS_2': ['hambre', 'alimentaciÁn', 'nutriciÁn'],
11147:             'ODS_3': ['salud', 'bienestar'],
11148:             'ODS_4': ['educaciÁn', 'calidad educativa'],
11149:             'ODS_5': ['igualdad de gÁnero', 'mujeres'],
11150:             'ODS_6': ['agua', 'saneamiento'],
11151:             'ODS_7': ['energÁa'],
11152:             'ODS_8': ['empleo', 'crecimiento econÁmico'],
11153:             'ODS_10': ['desigualdad', 'equidad'],
11154:             'ODS_11': ['ciudades sostenibles', 'desarrollo urbano'],
11155:             'ODS_13': ['cambio climÁtico', 'clima'],
11156:             'ODS_16': ['paz', 'justicia', 'instituciones']
11157:         }
11158:         text_lower = text.lower()
11159:         for sdg, themes in sdg_themes.items():
11160:             if any(theme in text_lower for theme in themes) and sdg not in sdgs:
11161:                 sdgs.append(sdg)
11162:
11163:     return sorted(list(set(sdgs)))
11164:
11165:     def _build_stakeholder_map(self, text: str) -> Dict[str, List[str]]:
11166:         """Construir mapa de stakeholders (simplificado)"""
11167:         stakeholders = defaultdict(list)
11168:         # Reutilizar inferencia de roles para llenar el mapa
11169:         entities = self._extract_policy_entities_with_context({'text': text})
11170:         for entity in entities:
11171:             role_name = entity.context_role.name.lower()
11172:             if entity.normalized_form not in stakeholders[role_name]:
11173:                 stakeholders[role_name].append(entity.normalized_form)
11174:         return dict(stakeholders)

```

```
11175:
11176:     def _calculate_policy_coherence(self, text: str, global_topics: Dict) -> float:
11177:         """Calcular coherencia de polÃtica (alineaciÃ³n con temas clave)"""
11178:         # SimplificaciÃ³n: media de la alineaciÃ³n tÃ³pica
11179:         return self._calculate_topic_alignment(text, global_topics)
11180:
11181:     def _extract_intervention_logic(self, strategic_unit: Dict) -> List[str]:
11182:         """Extraer la cadena lÃ³gica de intervenciÃ³n (e.g., insumo -> actividad -> producto -> resultado)"""
11183:         logic = []
11184:
11185:         # SimplificaciÃ³n: encadenamiento de antecedentes y consecuentes
11186:         for chain in strategic_unit.get('chains', []):
11187:             if chain.get('antecedent') and chain.get('consequent'):
11188:                 logic.append(f"{chain['antecedent']} -> {chain['consequent']}")
11189:
11190:         return logic[:5]
11191:
11192:     def _calculate_comprehensive_confidence(
11193:         self,
11194:         strategic_unit: Dict,
11195:         causal_evidence: List[CausalEvidence],
11196:         policy_entities: List[PolicyEntity]
11197:     ) -> Dict[str, float]:
11198:         """Calcular mÃ©tricas de confianza y calidad"""
11199:
11200:         weights = {
11201:             'causal': 0.4,
11202:             'entity': 0.3,
11203:             'structural': 0.2,
11204:             'semantic': 0.1
11205:         }
11206:
11207:         causal_conf = np.mean([e.confidence for e in causal_evidence]) if causal_evidence else 0.0
11208:         entity_conf = np.mean([e.confidence for e in policy_entities]) if policy_entities else 0.0
11209:         structural_conf = strategic_unit.get('confidence', 0.7)
11210:         semantic_conf = strategic_unit.get('semantic_coherence', 0.6)
11211:
11212:         overall = (
11213:             weights['causal'] * causal_conf +
11214:             weights['entity'] * entity_conf +
11215:             weights['structural'] * structural_conf +
11216:             weights['semantic'] * semantic_conf
11217:         )
11218:
11219:         return {
11220:             'causal_confidence': causal_conf,
11221:             'entity_confidence': entity_conf,
11222:             'structural_confidence': structural_conf,
11223:             'semantic_confidence': semantic_conf,
11224:             'overall_confidence': min(overall, 1.0)
11225:         }
11226:
11227:     def _find_cross_document_references(self, strategic_unit: Dict) -> List[CrossDocumentReference]:
11228:         """Encontrar referencias cruzadas documentales (Placeholder)"""
11229:         references = []
11230:         text = strategic_unit.get("text", "")
```

```

11231:     ref_patterns = [
11232:         (r'(?:ver|vÃ©ase)\s+(?:secciÃ³n|capÃ–tulo)\s+([^\.,\n]+)', 'explicit_reference'),
11233:         (r'como\s+se\s+(?:mencionÃ³|indicÃ³)\s+en\s+([^\.,\n]+)', 'backward_reference'),
11234:         (r'se\s+desarrollarÃ¡;\s+en\s+([^\.,\n]+)', 'forward_reference')
11235:     ]
11236:
11237:     for pattern, ref_type in ref_patterns:
11238:         matches = re.finditer(pattern, text, re.IGNORECASE)
11239:         for match in matches:
11240:             references.append(CrossDocumentReference(
11241:                 target_section=match.group(1).strip(),
11242:                 reference_type=ref_type,
11243:                 confidence=0.7,
11244:                 semantic_linkage=0.6,
11245:                 context_bridge=match.group(0)
11246:             ))
11247:
11248:     return references
11249:
11250: def _extract_implicit_assumptions(self, strategic_unit: Dict, causal_evidence: List[CausalEvidence]) -> List[Tuple[str, float]]:
11251:     """Extraer supuestos implÃ–citos (Placeholder)"""
11252:     assumptions = []
11253:     # Los antecedentes condicionales de baja confianza se consideran supuestos
11254:     for chain in strategic_unit.get('chains', []):
11255:         if chain.get('type') == 'conditional' and chain.get('confidence', 0.5) < 0.6:
11256:             assumptions.append((f"Asumiendo que: {chain['antecedent']}", 1.0 - chain.get('confidence', 0.5)))
11257:     return assumptions
11258:
11259: def _identify_contextual_presuppositions(self, strategic_unit: Dict) -> List[Tuple[str, float]]:
11260:     """Identificar presuposiciones contextuales (Placeholder)"""
11261:     presuppositions = []
11262:     text_lower = strategic_unit.get("text", "").lower()
11263:
11264:     if 'se continuarÃ¡;' in text_lower or 'marco existente' in text_lower:
11265:         presuppositions.append(("Existe un marco de polÃ–tica previo", 0.8))
11266:     if 'presupuesto asignado' in text_lower:
11267:         presuppositions.append(("Se cuenta con la disponibilidad de recursos", 0.9))
11268:
11269:     return presuppositions
11270:
11271: def _extract_discourse_markers(self, text: str) -> List[Tuple[str, str]]:
11272:     """Extraer marcadores de discurso (conectores)"""
11273:     markers = []
11274:
11275:     coherence_markers = {
11276:         'addition': [r'\by\b', r'\btambiÃ©n\b', r'\badicionalmente\b'],
11277:         'contrast': [r'\bpero\b', r'\bsin\s+embargo\b', r'\bno\s+obstante\b'],
11278:         'cause': [r'\bporque\b', r'\bya\s+que\b', r'\bdebido\s+a\b'],
11279:         'result': [r'\bpor\s+lo\s+tanto\b', r'\bpor\s+ende\b', r'\ben\s+consecuencia\b']
11280:     }
11281:
11282:     for relation_type, patterns in coherence_markers.items():
11283:         for pattern in patterns:
11284:             matches = re.finditer(pattern, text, re.IGNORECASE)
11285:             for match in matches:
11286:                 markers.append((match.group(0), relation_type))

```

```
11287:  
11288:         return markers  
11289:  
11290:     def _identify_rhetorical_patterns(self, text: str) -> List[str]:  
11291:         """Identificar patrones retóricos (e.g., afirmación, evidencia, conclusión)"""  
11292:         rhetorical = []  
11293:         rhetorical_patterns = self._analyze_rhetorical_structure(text)  
11294:  
11295:         for key, value in rhetorical_patterns.items():  
11296:             if value:  
11297:                 rhetorical.append(key)  
11298:  
11299:         return rhetorical  
11300:  
11301:     def _calculate_chunk_topic_distribution(self, text: str, global_topics: Dict) -> Dict[str, float]:  
11302:         """Calcular distribución temática del chunk"""  
11303:         distribution = {}  
11304:         text_lower = text.lower()  
11305:  
11306:         for topic in global_topics.get('topics', []):  
11307:             topic_score = 0  
11308:             for keyword, _ in topic['keywords'][:10]:  
11309:                 if keyword.lower() in text_lower:  
11310:                     topic_score += 1  
11311:             if topic_score > 0:  
11312:                 distribution[f"topic_{topic['topic_id']}"] = topic_score / 10.0  
11313:  
11314:         return distribution  
11315:  
11316:     def _extract_key_phrases(self, text: str) -> List[Tuple[str, float]]:  
11317:         """Extraer frases clave del texto (Placeholder: usando TF-IDF)"""  
11318:         try:  
11319:             # Usar TF-IDF para frases clave  
11320:             if not self.chunks_for_tfidf:  
11321:                 return []  
11322:  
11323:             tfidf_matrix = self.tfidf_vectorizer.fit_transform(self.chunks_for_tfidf)  
11324:             feature_names = self.tfidf_vectorizer.get_feature_names_out()  
11325:  
11326:             # Vectorizar solo el texto actual  
11327:             text_vector = self.tfidf_vectorizer.transform([text])  
11328:  
11329:             # Obtener los top features para este documento  
11330:             feature_array = text_vector.toarray().flatten()  
11331:             top_indices = feature_array.argsort()[-10:][::-1]  
11332:  
11333:             key_phrases = [(feature_names[i], feature_array[i]) for i in top_indices if feature_array[i] > 0]  
11334:             return key_phrases  
11335:  
11336:         except Exception as e:  
11337:             self.logger.error(f"Error en extracción de frases clave: {e}")  
11338:             return []  
11339:  
11340:     def _derive_kg_edges_for_chunk(self, entities: List[PolicyEntity], causal_evidence: List[CausalEvidence]) -> List[Tuple[str, str, str, float]]:  
11341:         """Derivar aristas del grafo de conocimiento para el chunk"""  
11342:         edges = []
```

```

11343:
11344:     # Entidades y relaciones causales
11345:     for entity1 in entities:
11346:         for entity2 in entities:
11347:             if entity1 != entity2:
11348:                 for evidence in causal_evidence:
11349:                     # Simple heuristic: si ambas entidades estÃ¡n en el contexto causal
11350:                     context = evidence.context_span
11351:                     if entity1.span[0] >= context[0] and entity2.span[1] <= context[1]:
11352:                         edges.append((
11353:                             entity1.normalized_form,
11354:                             entity2.normalized_form,
11355:                             evidence.causal_type.value,
11356:                             evidence.confidence
11357:                         ))
11358:
11359:     return edges[:10]
11360:
11361: def _get_model_versions(self) -> Dict[str, str]:
11362:     """Devuelve las versiones de los modelos utilizados"""
11363:     return {
11364:         'semantic_model': 'intfloat/multilingual-e5-large',
11365:         'spacy_model': self.nlp.meta.get('name') if self.nlp else 'None',
11366:         'pipeline_version': 'SMART-CHUNK-3.0-FINAL'
11367:     }
11368:
11369: def _calculate_semantic_density(self, text: str) -> float:
11370:     """Calcular la densidad semÃ¡ntica del chunk (placeholder)"""
11371:     # Densidad basada en la proporción de palabras clave
11372:     return min(len(re.findall(r'\b[A-Z][a-z]+\b', text)) / len(text.split()), 1.0)
11373:
11374: def _classify_chunk_type(self, text: str) -> ChunkType:
11375:     """Clasificar el tipo de chunk basado en el contenido (Placeholder)"""
11376:     text_lower = text.lower()
11377:     type_indicators = {
11378:         ChunkType.DIAGNOSTICO: ['diagnÃ³stico', 'anÃ¡lisis', 'situaciÃ³n actual'],
11379:         ChunkType.ESTRATEGIA: ['estrategia', 'objetivo', 'meta', 'propÃ³sito'],
11380:         ChunkType.METRICA: ['indicador', 'meta', 'mediciÃ³n', 'lÃnea base'],
11381:         ChunkType.FINANCIERO: ['presupuesto', 'recursos financieros', 'inversiÃ³n'],
11382:         ChunkType.NORMATIVO: ['ley', 'decreto', 'normativa', 'regulaciÃ³n'],
11383:         ChunkType.OPERATIVO: ['operaciÃ³n', 'implementaciÃ³n', 'ejecuciÃ³n'],
11384:         ChunkType.EVALUACION: ['evaluaciÃ³n', 'seguimiento', 'monitoreo']
11385:     }
11386:
11387:     scores = {}
11388:     for chunk_type, indicators in type_indicators.items():
11389:         score = sum(1 for ind in indicators if ind in text_lower)
11390:         if score > 0:
11391:             scores[chunk_type] = score
11392:
11393:     if scores:
11394:         return max(scores, key=scores.get)
11395:     return ChunkType.MIXTO
11396:
11397: def _calculate_coherence_score(self, strategic_unit: Dict) -> float:
11398:     """Calcular puntuaciÃ³n de coherencia"""

```

```

11399:         return strategic_unit.get('coherence', 0.75)
11400:
11401:     def _calculate_completeness_index(
11402:         self,
11403:         strategic_unit: Dict,
11404:         causal_evidence: List[CausalEvidence]
11405:     ) -> float:
11406:         """Calcular Índice de completitud"""
11407:         components = {
11408:             'has_text': bool(strategic_unit.get('text')),
11409:             'has_causal': len(causal_evidence) > 0,
11410:             'has_entities': len(strategic_unit.get('entities', [])) > 0,
11411:             'has_context': bool(strategic_unit.get('context')),
11412:             'has_hierarchy': bool(strategic_unit.get('hierarchy')),
11413:             'has_arg_structure': bool(strategic_unit.get('argument_structure')),
11414:             'has_temporal': bool(strategic_unit.get('temporal_dynamics'))
11415:         }
11416:
11417:         return sum(components.values()) / len(components)
11418:
11419: # --- MÁS de la clase principal (Continuación de smart_policy_chunks_industrial_v3_complete_final_Version2.py) ---
11420:
11421:     def _analyze_cross_references(self, text: str) -> Dict[str, List[Dict[str, Any]]]:
11422:         """Análisis de referencias cruzadas y relaciones inter-documento"""
11423:         references = defaultdict(list)
11424:
11425:         reference_patterns = [
11426:             (r'como\s+se\s+(?:mencion|indic|señal|estableci)\s+(?:en\s+)?(.+?)[,.,]', 'backward'),
11427:             (r'(?:(ver|v\@case|consultar)|\s+(?:la\s+)?(?:sección|capítulo|apartado))\s+(.+?)[,.,]', 'explicit'),
11428:             (r'tal\s+como\s+se\s+establece\s+en\s+(.+?)[,.,]', 'normative'),
11429:             (r'de\s+acuerdo\s+(?:con|a)\s+(?:lo\s+establecido\s+en\s+)?(.+?)[,.,]', 'normative'),
11430:             (r'según\s+(?:lo\s+dispuesto\s+en\s+)?(.+?)[,.,]', 'normative'),
11431:             (r'conforme\s+a\s+(.+?)[,.,]', 'normative'),
11432:             (r'en\s+el\s+marco\s+de\s+(.+?)[,.,]', 'framework'),
11433:             (r'se\s+desarrollará;\s+en\s+(.+?)[,.,]', 'forward')
11434:         ]
11435:
11436:         for pattern, ref_type in reference_patterns:
11437:             matches = re.finditer(pattern, text, re.IGNORECASE)
11438:             for match in matches:
11439:                 references[ref_type].append({
11440:                     'target': match.group(1).strip(),
11441:                     'position': match.span(),
11442:                     'context': text[max(0, match.start() - 100):min(len(text), match.end() + 100)]
11443:                 })
11444:
11445:         return dict(references)
11446:
11447:     def _analyze_temporal_structure(self, text: str) -> Dict[str, Any]:
11448:         """Análisis temporal completo"""
11449:         temporal_info = {
11450:             'time_markers': [],
11451:             'sequences': [],
11452:             'durations': [],
11453:             'milestones': [],
11454:             'temporal_ordering': []

```

```

11455:     }
11456:
11457:     # Marcadores temporales explÃ¡-citos
11458:     time_patterns = [
11459:         (r'\b(20\d{2})\b', 'year'),
11460:         (r'\b(enero|febrero|marzo|abril|mayo|junio|julio|agosto|septiembre|octubre|noviembre|diciembre)\s+(?:de\s+)?(20\d{2})\b', 'month_year'),
11461:         (r'\b(?:en|durante|hasta|desde|para)\s+(20\d{2})\b', 'temporal_prep'),
11462:         (r'\b(corto|mediano|largo)\s+plazo\b', 'horizon'),
11463:         (r'\b(trimestre|semestre|bimestre|cuatrimestre)\s+(\d+)\b', 'period'),
11464:         (r'\b(?:primer|segundo|tercer|cuarto)\s+(?:trimestre|semestre)\b', 'period_ordinal'),
11465:         (r'\b(inmediato|urgente|prioritario)\b', 'urgency'),
11466:         (r'\b(?:antes\s+de|despuÃ±@s\s+de|al\s+s+finalizar|una\s+vez\s+que)\b', 'sequence')
11467:     ]
11468:
11469:     for pattern, temp_type in time_patterns:
11470:         matches = re.finditer(pattern, text, re.IGNORECASE)
11471:         for match in matches:
11472:             temporal_info['time_markers'].append({
11473:                 'text': match.group(0),
11474:                 'type': temp_type,
11475:                 'position': match.span()
11476:             })
11477:
11478:     # Secuencias temporales
11479:     sequence_patterns = [
11480:         r'(?:(primero|en\s+primer\s+lugar)|',
11481:         r'(?:(segundo|en\s+segundo\s+lugar|luego|posteriormente)|',
11482:         r'(?:(tercero|en\s+tercer\s+lugar|finalmente|por\s+Ãºltimo)|'
11483:     ]
11484:
11485:     for idx, pattern in enumerate(sequence_patterns):
11486:         matches = re.finditer(pattern, text, re.IGNORECASE)
11487:         for match in matches:
11488:             temporal_info['sequences'].append({
11489:                 'order': idx + 1,
11490:                 'marker': match.group(0),
11491:                 'position': match.start()
11492:             })
11493:
11494:     return temporal_info
11495:
11496: def _analyze_discourse_structure(self, text: str) -> Dict[str, Any]:
11497:     """AnÃ¡lisis de estructura discursiva"""
11498:     discourse = {
11499:         'sentences': [],
11500:         'paragraphs': [],
11501:         'discourse_relations': [],
11502:         'rhetorical_moves': []
11503:     }
11504:
11505:     if self.nlp:
11506:         # Safe UTF-8 truncation for memory limit
11507:         truncated_text = safe_utf8_truncate(text, 500000)
11508:         doc = self.nlp(truncated_text)
11509:         discourse['sentences'] = [sent.text for sent in doc.sents][:1000]
11510:

```

```

11511:         # Extraer entidades nombradas
11512:         discourse['entities'] = [(ent.text, ent.label_) for ent in doc.ents][:200]
11513:
11514:         # Chunks nominales
11515:         discourse['noun_chunks'] = [chunk.text for chunk in doc.noun_chunks][:200]
11516:     else:
11517:         # Fallback sin SpaCy
11518:         discourse['sentences'] = re.split(r'[.!?]+\s+', text)[:1000]
11519:
11520:     # Detectar párrafos
11521:     paragraphs = text.split('\n\n')
11522:     discourse['paragraphs'] = [p.strip() for p in paragraphs if p.strip()[:500]
11523:
11524:     return discourse
11525:
11526: def _extract_coherence_relations(self, text: str) -> List[Dict[str, Any]]:
11527:     """Extraer relaciones de coherencia discursiva"""
11528:     relations = []
11529:
11530:     coherence_markers = {
11531:         'addition': [
11532:             r'\by\b', r'\be\b', r'\btambiÃ©n\b', r'\basimismo\b',
11533:             r'\badicionalmente\b', r'\ademÃ¡s\b', r'\bigualmente\b'
11534:         ],
11535:         'contrast': [
11536:             r'\bpero\b', r'\bsin\s+embargo\b', r'\bno\s+obstante\b',
11537:             r'\ben\s+cambio\b', r'\bpor\s+el\s+contrario\b', r'\bmientras\s+que\b'
11538:         ],
11539:         'cause': [
11540:             r'\bporque\b', r'\bya\s+que\b', r'\bdebido\s+a\b',
11541:             r'\bpuesto\s+que\b', r'\bdado\s+que\b', r'\ba\s+causa\s+de\b'
11542:         ],
11543:         'result': [
11544:             r'\bpor\s+lo\s+tanto\b', r'\bpor\s+ende\b', r'\basÃ–\b',
11545:             r'\ben\s+consecuencia\b', r'\bpor\s+coniguiente\b', r'\bde\s+modo\s+que\b'
11546:         ],
11547:         'condition': [
11548:             r'\bsi\b', r'\ben\s+caso\s+de\s+que\b', r'\bsiempre\s+que\b',
11549:             r'\bcon\s+tal\s+de\s+que\b', r'\ba\s+menos\s+que\b'
11550:         ],
11551:         'purpose': [
11552:             r'\bpara\s+que\b', r'\ba\s+fin\s+de\s+que\b', r'\bcon\s+el\s+objetivo\s+de\b',
11553:             r'\bcon\s+el\s+propÃ³sito\s+de\b'
11554:         ],
11555:         'elaboration': [
11556:             r'\bes\s+decir\b', r'\bo\s+sea\b', r'\ben\s+otras\s+palabras\b',
11557:             r'\bdicho\s+de\s+otro\s+modo\b'
11558:         ],
11559:         'example': [
11560:             r'\bpor\s+ejemplo\b', r'\bcomo\s+es\s+el\s+caso\s+de\b',
11561:             r'\btal\s+como\b', r'\bverbigracia\b'
11562:         ],
11563:     }
11564:
11565:     for relation_type, markers in coherence_markers.items():
11566:         for marker in markers:

```

```
11567:             matches = re.findall(marker, text, re.IGNORECASE)
11568:             for match in matches:
11569:                 relations.append({
11570:                     'type': relation_type,
11571:                     'marker': match.group(0),
11572:                     'position': match.span(),
11573:                     'confidence': 0.8
11574:                 })
11575:
11576:             return relations
11577:
11578:     def _analyze_rhetorical_structure(self, text: str) -> Dict[str, List[str]]:
11579:         """Analizar estructura retórica del texto"""
11580:         rhetorical = {
11581:             'claims': [],
11582:             'evidence': [],
11583:             'examples': [],
11584:             'conclusions': [],
11585:             'justifications': []
11586:         }
11587:
11588:         # Patrones de claims (afirmaciones)
11589:         claim_patterns = [
11590:             r'se\s+propone\s+que\s+[^.]+',
11591:             r'es\s+necesario\s+[^.]+',
11592:             r'se\s+debe\s+[^.]+',
11593:             r'resulta\s+fundamental\s+[^.]+',
11594:             r'se\s+requiere\s+[^.]+'
11595:         ]
11596:
11597:         for pattern in claim_patterns:
11598:             matches = re.findall(pattern, text, re.IGNORECASE)
11599:             rhetorical['claims'].extend(matches[:50])
11600:
11601:         # Patrones de evidencia
11602:         evidence_patterns = [
11603:             r'según\s+[^.]+',
11604:             r'de\s+acuerdo\s+con\s+[^.]+',
11605:             r'los\s+datos\s+(?:muestran|indican|demuestran)\s+[^.]+',
11606:             r'la\s+evidencia\s+(?:muestra|indica|demuestra)\s+[^.]+'
11607:         ]
11608:
11609:         for pattern in evidence_patterns:
11610:             matches = re.findall(pattern, text, re.IGNORECASE)
11611:             rhetorical['evidence'].extend(matches[:50])
11612:
11613:         # Patrones de ejemplos
11614:         example_patterns = [
11615:             r'por\s+ejemplo[^.]+',
11616:             r'como\s+es\s+el\s+caso\s+de\s+[^.]+',
11617:             r'tal\s+como\s+[^.]+'
11618:         ]
11619:
11620:         for pattern in example_patterns:
11621:             matches = re.findall(pattern, text, re.IGNORECASE)
11622:             rhetorical['examples'].extend(matches[:50])
```

```
11623:  
11624:     # Patrones de conclusiones  
11625:     conclusion_patterns = [  
11626:         r'en\s+conclusiÃ³n[^.]+',  
11627:         r'por\s+lo\s+tanto[^.]+',  
11628:         r'en\s+sÃ±tesis[^.]+'  
11629:     ]  
11630:  
11631:     for pattern in conclusion_patterns:  
11632:         matches = re.findall(pattern, text, re.IGNORECASE)  
11633:         rhetorical['conclusions'].extend(matches[:50])  
11634:  
11635:     return rhetorical  
11636:  
11637: def _analyze_information_flow(self, text: str) -> Dict[str, Any]:  
11638:     """  
11639:     Analyze information flow patterns in text.  
11640:  
11641:     Inputs:  
11642:         text (str): Text to analyze  
11643:     Outputs:  
11644:         Dict[str, Any]: Flow analysis metrics  
11645:     """  
11646:     sentences = filter_empty_sentences(re.split(r'[.!?]+\s+', text))  
11647:     sentences = [s for s in sentences if len(s) > 20]  
11648:     if not sentences:  
11649:         return {'sentence_count': 0}  
11650:  
11651:     sentence_lengths = [len(s.split()) for s in sentences]  
11652:     words = text.lower().split()  
11653:     unique_words = set(words)  
11654:  
11655:     return {  
11656:         'sentence_count': len(sentences),  
11657:         'avg_sentence_length': np.mean(sentence_lengths),  
11658:         'std_sentence_length': np.std(sentence_lengths),  
11659:         'lexical_diversity': len(unique_words) / len(words) if words else 0,  
11660:         'total_words': len(words),  
11661:         'unique_words': len(unique_words)  
11662:     } #  
11663:  
11664: def _extract_document_structure(self, text: str) -> Dict[str, Any]:  
11665:     """AnÃ¡lisis estructural del documento"""  
11666:     return {  
11667:         'section_hierarchy': self._identify_section_hierarchy(text),  
11668:         'policy_frameworks': self._identify_policy_frameworks(text),  
11669:         'raw_text': text # Almacenar texto para referencia  
11670:     }  
11671:  
11672: def _identify_section_hierarchy(self, text: str) -> List[Dict[str, Any]]:  
11673:     """Identificar la jerarquÃ-a de secciones/tÃ-tulos"""  
11674:     hierarchy = []  
11675:     # Patrones de encabezados (simplificados)  
11676:     patterns = [  
11677:         (r'^(\CAPÃ\u215TULO\s+I[VX]+\s*:\s*.+)\$', 1, 'chapter'),  
11678:         (r'^(\SECCIÃ\u2153N\s+[A-Z]\s*:\s*.+)\$', 2, 'section'),
```

```

11679:             (r'^\d+\.\d+(?:\.\d+)?s+([A-ZÀ\221Ã\201Ã\211Ã\215Ã\223Ã\232][^.!?]+'\$, 3, 'header_numbered'),
11680:             (r'^([a-z])\)\s+([A-ZÀ\221Ã\201Ã\211Ã\215Ã\223Ã\232][^.!?]+'\$, 4, 'item_alpha'),
11681:             (r'^à\200\$s+([A-ZÀ\221Ã\201Ã\211Ã\215Ã\223Ã\232][^.!?]+'\$'\$, 5, 'bullet'),
11682:         ]
11683:
11684:     lines = text.split('\n')
11685:     for idx, line in enumerate(lines):
11686:         line = line.strip()
11687:         if not line or len(line) < 3: continue
11688:
11689:         for pattern, level, header_type in patterns:
11690:             match = re.match(pattern, line, re.MULTILINE)
11691:             if match:
11692:                 title = match.group(2) if match.lastindex >= 2 else match.group(1)
11693:                 hierarchy.append({
11694:                     'title': title.strip(),
11695:                     'level': level,
11696:                     'line_number': idx,
11697:                     'type': header_type,
11698:                     'full_text': line
11699:                 })
11700:             break
11701:     return hierarchy
11702:
11703: def _identify_policy_frameworks(self, text: str) -> List[Dict[str, Any]]:
11704:     """Identificar todos los marcos de polÃtica presentes"""
11705:     frameworks = []
11706:     framework_patterns = [
11707:         (r'plan\s+(:de\s+)?desarrollo\s+(:municipal|departamental|nacional)?', 'plan_desarrollo'),
11708:         (r'polÃtica\s+pÃblica\s+(:de\s+)?[\w\s]+', 'politica_publica'),
11709:         (r'ley\s+[\w\s]+', 'legal_framework')
11710:     ]
11711:
11712:     for pattern, framework_type in framework_patterns:
11713:         matches = re.finditer(pattern, text, re.IGNORECASE)
11714:         for match in matches:
11715:             frameworks.append({
11716:                 'text': match.group(0),
11717:                 'type': framework_type,
11718:                 'position': match.span()
11719:             })
11720:     return frameworks
11721:
11722: def _extract_content_for_pa_dimension(
11723:     self,
11724:     document_text: str,
11725:     policy_area: str,
11726:     dimension: str,
11727:     sentences: List[str],
11728:     sentence_positions: List[Tuple[int, int]],
11729:     sentence_embeddings: Optional[np.ndarray] = None
11730: ) -> Dict[str, Any]:
11731:     """
11732:     Extract most relevant content for a specific (PA, DIM) combination.
11733:
11734:     Uses embedding-based similarity to find sentences most aligned with

```

```
11735:     the policy area and dimension keywords.
11736:
11737:     Args:
11738:         document_text: Full document text
11739:         policy_area: Policy area code (PA01-PA10)
11740:         dimension: Dimension code (DIM01-DIM06)
11741:         sentences: List of document sentences
11742:         sentence_positions: List of (start, end) byte positions for each sentence
11743:         sentence_embeddings: Pre-computed embeddings (optional, computed if None)
11744:
11745:     Returns:
11746:         Dictionary with segment metadata and text
11747:     """
11748:     # Generate query embedding from PA + DIM keywords
11749:     pa_keywords = self._pa_keywords.get(policy_area, [])
11750:     dim_keywords = self._dim_keywords.get(dimension, [])
11751:     query_text = " ".join(pa_keywords + dim_keywords)
11752:
11753:     # Get query embedding
11754:     query_embedding = self._spc_sem.embed_text(query_text)
11755:
11756:     # Compute sentence embeddings if not provided
11757:     if sentence_embeddings is None:
11758:         sentence_embeddings = self._spc_sem.embed_batch(sentences)
11759:
11760:     # Compute cosine similarity between query and all sentences
11761:     from sklearn.metrics.pairwise import cosine_similarity
11762:     similarities = cosine_similarity(
11763:         query_embedding.reshape(1, -1),
11764:         sentence_embeddings
11765:     )[0]
11766:
11767:     # Find top-K most similar sentences (K=10)
11768:     top_k = min(10, len(sentences))
11769:     top_indices = np.argsort(similarities)[-top_k:][::-1]
11770:
11771:     # Extract contiguous region around top sentences
11772:     if len(top_indices) == 0:
11773:         # Fallback: return first 800 chars
11774:         segment_text = document_text[:800]
11775:         segment_start = 0
11776:         segment_end = len(segment_text)
11777:     else:
11778:         # Find min/max positions to create contiguous chunk
11779:         min_idx = min(top_indices)
11780:         max_idx = max(top_indices)
11781:
11782:         # Expand window by ±2 sentences for context
11783:         start_idx = max(0, min_idx - 2)
11784:         end_idx = min(len(sentences)) - 1, max_idx + 2
11785:
11786:         # Get byte positions
11787:         segment_start = sentence_positions[start_idx][0]
11788:         segment_end = sentence_positions[end_idx][1]
11789:         segment_text = document_text[segment_start:segment_end]
11790:
```

```

11791:     # Compute relevance score (mean of top-K similarities)
11792:     relevance_score = float(np.mean(similarities[top_indices])) if len(top_indices) > 0 else 0.0
11793:
11794:     return {
11795:         "text": segment_text,
11796:         "position": (segment_start, segment_end),
11797:         "policy_area": policy_area,
11798:         "dimension": dimension,
11799:         "relevance_score": relevance_score,
11800:         "top_sentence_indices": top_indices.tolist(),
11801:         "query_keywords": pa_keywords + dim_keywords
11802:     }
11803:
11804: def _generate_60_structured_segments(
11805:     self,
11806:     document_text: str,
11807:     structural_analysis: Dict[str, Any]
11808: ) -> List[Dict[str, Any]]:
11809: """
11810:     Generate EXACTLY 60 structured segments aligned by (PA \ 227 DIM) matrix.
11811:
11812:     This replaces FASE 4 semantic segmentation with structured extraction.
11813:     Each segment is explicitly aligned to one Policy Area and one Dimension.
11814:
11815:     Args:
11816:         document_text: Full document text
11817:         structural_analysis: Output from FASE 3 (document structure analysis)
11818:
11819:     Returns:
11820:         List of 60 segment dictionaries, one per (PA, DIM) combination
11821: """
11822: from farfan_pipeline.core.canonical_notation import get_all_policy_areas, get_all_dimensions
11823:
11824: self.logger.info("FASE 4: Generating 60 structured segments (PA \ 227 DIM matrix)")
11825:
11826: # Split document into sentences for embedding-based extraction
11827: sentences = sent_tokenize(document_text, language='spanish')
11828:
11829: # Compute sentence positions
11830: sentence_positions = []
11831: current_pos = 0
11832: for sent in sentences:
11833:     start = document_text.find(sent, current_pos)
11834:     end = start + len(sent)
11835:     sentence_positions.append((start, end))
11836:     current_pos = end
11837:
11838: # Pre-compute sentence embeddings (once for all 60 extractions)
11839: self.logger.info(f"Computing embeddings for {len(sentences)} sentences...")
11840: sentence_embeddings = self._spc_sem.embed_batch(sentences)
11841:
11842: # Generate 60 segments
11843: policy_areas = get_all_policy_areas() # PA01..PA10
11844: dimensions = get_all_dimensions() # D1..D6
11845:
11846: structured_segments = []

```

```
11847:  
11848:     for pa_code, pa_info in policy_areas.items():  
11849:         for dim_key, dim_info in dimensions.items():  
11850:             segment = self._extract_content_for_pa_dimension(  
11851:                 document_text=document_text,  
11852:                 policy_area=pa_code,  
11853:                 dimension=dim_info.code,  
11854:                 sentences=sentences,  
11855:                 sentence_positions=sentence_positions,  
11856:                 sentence_embeddings=sentence_embeddings  
11857:             )  
11858:  
11859:             # Add PA and DIM metadata  
11860:             segment['policy_area_id'] = pa_code  
11861:             segment['dimension_id'] = dim_info.code  
11862:             segment['pa_name'] = pa_info.name  
11863:             segment['dim_label'] = dim_info.label  
11864:  
11865:             structured_segments.append(segment)  
11866:  
11867:             self.logger.debug(  
11868:                 f" Extracted segment for {pa_code} \u2277 {dim_info.code}: "  
11869:                 f"{len(segment['text'])} chars, relevance={segment['relevance_score']:.3f}"  
11870:             )  
11871:  
11872:             assert len(structured_segments) == 60, \  
11873:                 f"Expected 60 segments, got {len(structured_segments)}"  
11874:  
11875:             self.logger.info(f"\u234\u205 Generated exactly {len(structured_segments)} structured segments")  
11876:  
11877:             return structured_segments  
11878:  
11879: def generate_smart_chunks(self, document_text: str, document_metadata: Dict) -> List[SmartPolicyChunk]:  
11880:     """  
11881:     Main pipeline phase: Generate 60 Smart Policy Chunks aligned with the P01-ES v1.0 spec.  
11882:  
11883:     This streamlined pipeline executes the structured (Policy Area x Dimension)  
11884:     segmentation and bypasses the subsequent complex analysis to ensure the  
11885:     output is exactly the 60 thematically isolated chunks required by the orchestrator.  
11886:     """  
11887:     self.logger.info("Starting Smart Chunks v3.0 pipeline (P01-ES v1.0 Compliant Mode)")  
11888:  
11889:     # FASE 1: Preprocesamiento avanzado  
11890:     normalized_text = self._advanced_preprocessing(document_text)  
11891:  
11892:     # FASE 2: Análisis estructural y de jerarquía (needed for segmentation context)  
11893:     structural_analysis = self._extract_document_structure(normalized_text)  
11894:     structural_analysis['raw_text'] = document_text  
11895:  
11896:     # FASE 4: Structured (PA \u2277 DIM) segmentation - EXACTLY 60 chunks  
11897:     structured_segments = self._generate_60_structured_segments(  
11898:         document_text=document_text,  
11899:         structural_analysis=structural_analysis  
11900:     )  
11901:     self.logger.info(f"\u234\u205 Generated {len(structured_segments)} structured segments (PA \u2277 DIM)")  
11902:
```

```

11903:     # Convert the 60 segments to SmartPolicyChunk objects
11904:     smart_chunks = []
11905:     document_id = document_metadata.get("document_id", "doc_001")
11906:
11907:     for segment in structured_segments:
11908:         text = segment.get("text", "")
11909:         content_hash = hashlib.sha256(text.encode('utf-8')).hexdigest()
11910:         pa_id = segment.get("policy_area_id")
11911:         dim_id = segment.get("dimension_id")
11912:         chunk_id = f"{document_id}_{pa_id}_{dim_id}_{content_hash[:8]}"
11913:
11914:         chunk = SmartPolicyChunk(
11915:             chunk_id=chunk_id,
11916:             document_id=document_id,
11917:             content_hash=content_hash,
11918:             text=text,
11919:             normalized_text=self._advanced_preprocessing(text),
11920:             policy_area_id=pa_id,
11921:             dimension_id=dim_id,
11922:             document_position=segment.get("position", (0, 0)),
11923:             # --- Default empty values for bypassed analysis ---
11924:             semantic_density=0.0,
11925:             section_hierarchy=[],
11926:             chunk_type=ChunkType.MIXTO,
11927:             causal_chain=[],
11928:             policy_entities=[],
11929:             implicit_assumptions=[],
11930:             contextual_presuppositions=[],
11931:             confidence_metrics={'overall_confidence': segment.get('relevance_score', 0.0)},
11932:             coherence_score=segment.get('relevance_score', 0.0),
11933:             strategic_importance=segment.get('relevance_score', 0.0),
11934:         )
11935:         smart_chunks.append(chunk)
11936:
11937:     self.logger.info(f"Pipeline completed: {len(smart_chunks)} chunks generated.")
11938:     return smart_chunks
11939:
11940: def _advanced_preprocessing(self, text: str) -> str:
11941:     """
11942:     Advanced text preprocessing with normalization and encoding fixes.
11943:
11944:     Inputs:
11945:         text (str): Raw input text
11946:     Outputs:
11947:         str: Normalized and cleaned text
11948:
11949:     # NormalizaciÃ³n de espacios y saltos de lÃnea
11950:     text = re.sub(r'\s+', ' ', text)
11951:     text = re.sub(r'\n+', '\n', text)
11952:
11953:     # CorrecciÃ³n de encodings problemÃ¡ticos (common UTF-8 mojibake)
11954:     encoding_fixes = {
11955:         '\u203a': '\u00e1', '\u203c': '\u00e0', '\u203d': '\u00e2', '\u203b': '\u00e3', '\u203e': '\u00e0',
11956:         '\u203f': '\u00e1', '\u2039': '\u00e2', '\u2038': '\u00e3', '\u2037': '\u00e0', '\u2036': '\u00e2',
11957:         '\u2035': '\u00e1', '\u2034': '\u00e2', '\u2033': '\u00e3', '\u2032': '\u00e0', '\u2031': '\u00e1',
11958:         '\u2030': '\u00e2', '\u2032': '\u00e3', '\u2031': '\u00e0', '\u2033': '\u00e1'
11959:     }

```

```

11959:         }
11960:         for wrong, correct in encoding_fixes.items():
11961:             text = text.replace(wrong, correct)
11962:
11963:         return text
11964:
11965:     def _enrich_with_inter_chunk_relationships(self, chunks: List[SmartPolicyChunk]) -> List[SmartPolicyChunk]:
11966:         """
11967:             CANONICAL SOTA: Enrich chunks with semantic relationships using BGE-M3 embeddings.
11968:
11969:             Uses batch embeddings and vectorized similarity computation instead of
11970:             manual loops and sklearn.cosine_similarity.
11971:             """
11972:         if not chunks:
11973:             return []
11974:
11975:         texts = [c.text for c in chunks]
11976:         self.corpus_embeddings = self._generate_embeddings_for_corpus(texts)
11977:
11978:         # Vectorized cosine similarity (no sklearn dependency)
11979:         norms = np.linalg.norm(self.corpus_embeddings, axis=1, keepdims=True)
11980:         normalized_embs = self.corpus_embeddings / (norms + 1e-8)
11981:         similarity_matrix = np.dot(normalized_embs, normalized_embs.T)
11982:
11983:         # Efficiently find related chunks using vectorized operations
11984:         for i, chunk in enumerate(chunks):
11985:             # Get similarities for this chunk (excluding self)
11986:             sims = similarity_matrix[i].copy()
11987:             sims[i] = -1 # Exclude self
11988:
11989:             # Find chunks above threshold
11990:             related_indices = np.where(sims >= self.config.CROSS_REFERENCE_MIN_SIMILARITY)[0]
11991:
11992:             # Sort by similarity
11993:             sorted_indices = related_indices[np.argsort(-sims[related_indices])]
11994:
11995:             chunk.related_chunks = [
11996:                 (chunks[j].chunk_id, float(sims[j]))
11997:                 for j in sorted_indices
11998:             ]
11999:
12000:         return chunks
12001:
12002:     def _generate_embeddings_for_corpus(self, texts: List[str], batch_size: int = 64) -> np.ndarray:
12003:         """
12004:             CANONICAL SOTA: Batch embeddings using SemanticChunkingProducer.
12005:
12006:             SemanticChunkingProducer handles batching internally with BGE-M3.
12007:             batch_size kept for signature compatibility.
12008:
12009:             Inputs:
12010:                 texts (List[str]): List of text strings to embed
12011:                 batch_size (int): Batch size hint (default: 64)
12012:             Outputs:
12013:                 np.ndarray: Array of embeddings, shape (n_texts, embedding_dim)
12014:             """

```

```

12015:         if not texts:
12016:             return np.array([])
12017:
12018:         # CANONICAL SOTA: Use SemanticChunkingProducer batch embedding
12019:         embs = self._spc_sem.embed_batch(texts)
12020:         return np.vstack(embs).astype(np.float32)
12021:
12022:     def _validate_strategic_integrity(self, chunks: List[SmartPolicyChunk]) -> List[SmartPolicyChunk]:
12023:         """Validar que los chunks cumplan con umbrales mÁ-nimos de calidad y completitud"""
12024:         validated = []
12025:         for chunk in chunks:
12026:             # Criterios de validaciÃ³n
12027:             is_valid = (
12028:                 chunk.coherence_score >= self.config.MIN_COHERENCE_SCORE and
12029:                 chunk.completeness_index >= self.config.MIN_COMPLETENESS_INDEX and
12030:                 chunk.strategic_importance >= self.config.MIN_STRATEGIC_IMPORTANCE and
12031:                 chunk.information_density >= self.config.MIN_INFORMATION_DENSITY and
12032:                 len(chunk.text) >= self.config.MIN_CHUNK_SIZE
12033:             )
12034:
12035:             if is_valid:
12036:                 validated.append(chunk)
12037:             else:
12038:                 self.logger.warning(f"Chunk {chunk.chunk_id} no pasÃ³ validaciÃ³n de integridad")
12039:
12040:         return validated
12041:
12042:     def _intelligent_deduplication(self, chunks: List[SmartPolicyChunk]) -> List[SmartPolicyChunk]:
12043:         """DeduplicaciÃ³n inteligente de chunks"""
12044:         if len(chunks) < 2:
12045:             return chunks
12046:
12047:         deduplicated = []
12048:         processed_hashes = set()
12049:
12050:         for chunk in chunks:
12051:             # Verificar por hash exacto
12052:             if chunk.content_hash not in processed_hashes:
12053:                 processed_hashes.add(chunk.content_hash)
12054:                 deduplicated.append(chunk)
12055:
12056:         # DeduplicaciÃ³n semÃ¡ntica (para near-duplicates) using vectorized ops
12057:         final_list = []
12058:         if deduplicated:
12059:             embeddings = self._generate_embeddings_for_corpus([c.text for c in deduplicated])
12060:             n_chunks = len(deduplicated)
12061:
12062:             # Vectorized cosine similarity (no sklearn dependency)
12063:             norms = np.linalg.norm(embeddings, axis=1, keepdims=True)
12064:             normalized_embs = embeddings / (norms + 1e-8)
12065:             sim_matrix = np.dot(normalized_embs, normalized_embs.T)
12066:
12067:             is_duplicate = [False] * n_chunks
12068:
12069:             for i in range(n_chunks):
12070:                 if is_duplicate[i]:

```

```

12071:         continue
12072:
12073:         final_list.append(deduplicated[i])
12074:
12075:         for j in range(i + 1, n_chunks):
12076:             if sim_matrix[i, j] >= self.config.DEDUPLICATION_THRESHOLD:
12077:                 is_duplicate[j] = True
12078:
12079:     return final_list
12080:
12081:     def _rank_by_strategic_importance(self, chunks: List[SmartPolicyChunk]) -> List[SmartPolicyChunk]:
12082:         """Ranking final de chunks por importancia estratÃ©gica"""
12083:         # Usar el campo 'strategic_importance' calculado en _create_smart_policy_chunk
12084:         chunks.sort(key=lambda x: x.strategic_importance, reverse=True)
12085:         return chunks
12086:
12087:     # --- MÃ©todos de EvaluaciÃ³n (ContinuaciÃ³n de smart_policy_chunks_industrial_v3_part4_completion.py) ---
12088:
12089:     def _assess_strategic_importance(
12090:         self,
12091:         strategic_unit: Dict,
12092:         causal_evidence: List[CausalEvidence],
12093:         policy_entities: List[PolicyEntity]
12094:     ) -> float:
12095:         """Evaluar importancia estratÃ©gica"""
12096:         factors = {
12097:             'causal_strength': np.mean([e.strength_score for e in causal_evidence]) if causal_evidence else 0,
12098:             'entity_relevance': len([e for e in policy_entities if e.context_role in [PolicyEntityRole.EXECUTOR, PolicyEntityRole.BENEFICIARY]]) / max(len(policy_entities), 1),
12099:             'position_weight': 1.0 if strategic_unit.get('position', (0, 0))[0] < 1000 else 0.5, # Ponderar texto temprano
12100:             'keyword_presence': self._calculate_strategic_keyword_presence(strategic_unit.get('text', ''))
12101:         }
12102:
12103:         return np.mean(list(factors.values()))
12104:
12105:     def _calculate_strategic_keyword_presence(self, text: str) -> float:
12106:         """Calcular presencia de palabras clave estratÃ©gicas"""
12107:         strategic_keywords = [
12108:             'objetivo', 'meta', 'estrategia', 'prioridad', 'desarrollo',
12109:             'transformaciÃ³n', 'impacto', 'resultado', 'sostenible', 'integral'
12110:         ]
12111:
12112:         text_lower = text.lower()
12113:         keyword_count = sum(1 for kw in strategic_keywords if kw in text_lower)
12114:
12115:         return min(keyword_count / len(strategic_keywords), 1.0)
12116:
12117:     def _calculate_information_density(self, text: str) -> float:
12118:         """Calcular densidad de informaciÃ³n"""
12119:         if not text:
12120:             return 0.0
12121:
12122:         # MÃ©tricas de densidad
12123:         words = text.split()
12124:         sentences = re.split(r'[.!?]+', text)
12125:

```

```
12126:         if not words or not sentences:
12127:             return 0.0
12128:
12129:         metrics = {
12130:             'lexical_diversity': len(set(words)) / len(words),
12131:             'avg_sentence_length': np.mean([len(s.split()) for s in sentences if s.strip()]),
12132:             'entity_density': len(re.findall(r'\b[A-Z][a-z]+\b', text)) / len(words) # EstimaciÃ³n de entidades
12133:         }
12134:
12135:         return np.mean(list(metrics.values()))
12136:
12137:     def _assess_actionability(self, text: str, policy_entities: List[PolicyEntity]) -> float:
12138:         """Evaluar accionabilidad del chunk"""
12139:         action_indicators = [
12140:             r'se\s+(?:debe|deberÃ;|requiere)',
12141:             r'es\s+necesario',
12142:             r'(?:(?:implementar|ejecutar|desarrollar|crear|establecer|asignar)\s+',
12143:             r'(?:(?:el|la)\s+(?:Ministerio|AlcaldÃ-a|DirecciÃ³n)\s+(?:debe|deberÃ;))'
12144:         ]
12145:
12146:         action_score = sum(len(re.findall(pat, text, re.IGNORECASE)) for pat in action_indicators)
12147:
12148:         # Presencia de ejecutores explÃ–citos
12149:         executor_count = len([e for e in policy_entities if e.context_role == PolicyEntityRole.EXECUTOR])
12150:
12151:         final_score = (min(action_score / 5.0, 1.0) * 0.7) + (min(executor_count / 2.0, 1.0) * 0.3)
12152:         return final_score
12153:
12154: # =====
12155: # CANONICAL PIPELINE WRAPPER
12156: # =====
12157:
12158: class CPPIngestionPipeline:
12159:     """
12160:         Canonical wrapper for the SPC ingestion process.
12161:         Provides the interface expected by run_policy_pipeline_verified.py.
12162:     """
12163:     def __init__(self):
12164:         self.system = StrategicChunkingSystem()
12165:         self.logger = logging.getLogger("CPPIngestionPipeline")
12166:
12167:     @async def process(self, pdf_path: Path) -> Any:
12168:         """
12169:             Process a PDF document into a CanonPolicyPackage (SPC output).
12170:
12171:             Args:
12172:                 pdf_path: Path to the input PDF file.
12173:
12174:             Returns:
12175:                 The processed CanonPolicyPackage (or equivalent SPC object).
12176:         """
12177:         self.logger.info(f"Starting CPPIngestionPipeline for: {pdf_path}")
12178:
12179:         if not pdf_path.exists():
12180:             raise FileNotFoundError(f"Input PDF not found: {pdf_path}")
12181:
```

```
12182:     # Extract text with provenance (page numbers)
12183:     # Using pdfplumber as per SOTA requirements
12184:     document_text = ""
12185:     page_map = [] # List of (start_char, end_char, page_number)
12186:
12187:     try:
12188:         import pdfplumber
12189:         with pdfplumber.open(pdf_path) as pdf:
12190:             current_char = 0
12191:             for i, page in enumerate(pdf.pages):
12192:                 text = page.extract_text() or ""
12193:                 # Normalize text slightly to avoid issues
12194:                 text = self.system._advanced_preprocessing(text)
12195:
12196:                 start = current_char
12197:                 end = start + len(text)
12198:                 page_number = i + 1
12199:
12200:                 page_map.append({
12201:                     "start": start,
12202:                     "end": end,
12203:                     "page": page_number,
12204:                     "bbox": page.bbox
12205:                 })
12206:
12207:                 document_text += text + "\n" # Add newline as separator
12208:                 current_char = len(document_text)
12209:
12210:     except ImportError:
12211:         self.logger.warning("pdfplumber not installed. Fallback to basic text reading (NO PROVENANCE).")
12212:         # Fallback for dev/testing without pdfplumber
12213:         try:
12214:             document_text = pdf_path.read_text(encoding='utf-8')
12215:         except UnicodeDecodeError:
12216:             # Try latin-1 as fallback
12217:             document_text = pdf_path.read_text(encoding='latin-1')
12218:
12219:     # Generate chunks using the core system
12220:     # Note: We pass the full text. The system currently doesn't natively handle
12221:     # the page_map during chunk generation, so we need to inject provenance *after*
12222:     # chunking but *before* returning.
12223:
12224:     # Create a mock metadata dict
12225:     metadata = {
12226:         "source_file": str(pdf_path),
12227:         "ingestion_timestamp": canonical_timestamp()
12228:     }
12229:
12230:     # Run the core pipeline
12231:     chunks = self.system.generate_smart_chunks(document_text, metadata)
12232:
12233:     # POST-PROCESSING: Inject Provenance
12234:     # Map chunks back to pages using the page_map
12235:     if page_map:
12236:         self._inject_provenance(chunks, page_map, str(pdf_path))
12237:     else:
```



```
12294:             break
12295:
12296:         if primary_page:
12297:             chunk.provenance = Provenance(
12298:                 page_number=primary_page,
12299:                 section_header=None, # Could be extracted if we had section info
12300:                 source_file=source_file
12301:             )
12302:
12303: # =====
12304: # SCRIPT DE EJECUCIÃ“N (MAIN)
12305: # =====
12306:
12307: def validate_cli_arguments(args):
12308:     """
12309:     Validate CLI arguments before processing.
12310:
12311:     Returns:
12312:         Tuple[Path, Path]: Normalized input and output paths
12313:     """
12314:     input_path = Path(args.input).expanduser()
12315:     if not input_path.exists():
12316:         raise ValidationError(f"Input file not found: {input_path}")
12317:
12318:     output_path = Path(args.output).expanduser()
12319:     output_dir = output_path.parent if output_path.name else output_path
12320:     if not output_dir.exists():
12321:         raise ValidationError(f"Output directory does not exist: {output_dir}")
12322:     if not os.access(output_dir, os.W_OK):
12323:         raise ValidationError(f"Output directory is not writable: {output_dir}")
12324:
12325:     if not args.doc_id or not args.doc_id.strip():
12326:         raise ValidationError("Document ID cannot be empty")
12327:
12328:     if args.max_chunks < 0:
12329:         raise ValidationError(f"max_chunks must be non-negative, got: {args.max_chunks}")
12330:
12331:     logger.info("CLI arguments validated successfully")
12332:     return input_path, output_path
12333:
12334: def main(args):
12335:     """
12336:     Main pipeline function for Smart Policy Chunks generation.
12337:
12338:     Inputs:
12339:         args: Command-line arguments with input/output paths and configuration
12340:     Outputs:
12341:         int: Exit code (0 for success, 1 for error)
12342:     """
12343:     try:
12344:         # 0. Validate CLI arguments
12345:         input_path, output_path = validate_cli_arguments(args)
12346:
12347:         # 1. Cargar el documento
12348:         logger.info(f"Loading input document: {input_path}")
12349:         try:
```

```
12350:     document_text = input_path.read_text(encoding='utf-8')
12351: except IOError as e:
12352:     raise ProcessingError(f"Failed to read input file: {e}")
12353:
12354: # 2. Inicializar el sistema (con lazy loading)
12355: logger.info("Initializing Strategic Chunking System...")
12356: chunking_system = StrategicChunkingSystem()
12357:
12358: # Metadata del documento
12359: metadata = {
12360:     'document_id': args.doc_id,
12361:     'title': args.title,
12362:     'version': 'v3.0',
12363:     'processing_timestamp': canonical_timestamp()
12364: }
12365:
12366: # 3. Generar Smart Chunks
12367: logger.info("Generating smart policy chunks...")
12368: chunks = chunking_system.generate_smart_chunks(document_text, metadata)
12369:
12370: # 4. Limitar el nÃºmero de chunks a guardar si se especifica
12371: if args.max_chunks > 0:
12372:     logger.info(f"Limiting output to {args.max_chunks} chunks")
12373:     chunks = chunks[:args.max_chunks]
12374:
12375: # 5. Serializar resultados (summary version with truncated text)
12376: logger.info(f"Generated {len(chunks)} chunks. Saving results...")
12377: output_data = {
12378:     'metadata': metadata,
12379:     'config': {
12380:         'min_chunk_size': chunking_system.config.MIN_CHUNK_SIZE,
12381:         'max_chunk_size': chunking_system.config.MAX_CHUNK_SIZE,
12382:         'semantic_coherence_threshold': chunking_system.config.SEMANTIC_COHERENCE_THRESHOLD
12383:     },
12384:     'chunks': [
12385:         asdict(c) for c in chunks
12386:     ]
12387: }
12388:
12389: # Truncar el texto del chunk de forma segura para summary.json
12390: for chunk_data in output_data['chunks']:
12391:     original_text = chunk_data['text']
12392:     # Safe UTF-8 truncation to 200 bytes
12393:     chunk_data['text'] = safe_utf8_truncate(original_text, 200)
12394:     if len(original_text.encode('utf-8')) > 200:
12395:         chunk_data['text'] += '...'
12396:
12397: summary_payload = json.dumps(output_data, indent=2, ensure_ascii=False, default=np_to_list)
12398: try:
12399:     output_path.write_text(summary_payload, encoding='utf-8')
12400:     logger.info(f"Summary output saved to: {output_path}")
12401: except (IOError, TypeError) as e:
12402:     raise SerializationError(f"Failed to write summary output: {e}")
12403:
12404: # 6. Guardar versiÃ³n completa sin truncar texto (full.json)
12405: if args.save_full:
```

```
12406:         full_output = output_path.with_name(f"{output_path.stem}_full.json")
12407:     try:
12408:         # Crear una estructura completa con texto Ántegro
12409:         full_data = {
12410:             'metadata': metadata,
12411:             'config': output_data['config'],
12412:             'chunks': []
12413:         }
12414:
12415:         for chunk in chunks:
12416:             chunk_dict = asdict(chunk)
12417:             # Mantener texto completo
12418:             full_data['chunks'].append(chunk_dict)
12419:
12420:             full_payload = json.dumps(full_data, indent=2, ensure_ascii=False, default=np_to_list)
12421:             full_output.write_text(full_payload, encoding='utf-8')
12422:             logger.info(f"Full output saved to: {full_output}")
12423:         except (IOError, TypeError) as e:
12424:             raise SerializationError(f"Failed to write full output: {e}")
12425:
12426: # 7. Generar reporte de verificación con canonical timestamps
12427: verification = {
12428:     'pipeline_version': 'SMART-CHUNK-3.0-FINAL',
12429:     'execution_timestamp': canonical_timestamp(),
12430:     'input_file': args.input,
12431:     'input_hash': hashlib.sha256(document_text.encode('utf-8')).hexdigest(),
12432:     'output_hash': hashlib.sha256(summary_payload.encode('utf-8')).hexdigest(),
12433:     'chunks_generated': len(chunks),
12434:     'validation_passed': all(
12435:         c.coherence_score >= chunking_system.config.MIN_COHERENCE_SCORE
12436:         for c in chunks
12437:     ),
12438:     'success': len(chunks) > 0 and all(c.coherence_score >= chunking_system.config.MIN_COHERENCE_SCORE for c in chunks)
12439: }
12440:
12441: verification_file = output_path.with_name(f"{output_path.stem}_verification.json")
12442: verification_payload = json.dumps(verification, indent=2, default=np_to_list)
12443: try:
12444:     verification_file.write_text(verification_payload, encoding='utf-8')
12445:     logger.info(f"Verification report saved to: {verification_file}")
12446: except IOError as e:
12447:     logger.warning(f"Failed to write verification file: {e}")
12448:
12449: if verification['success']:
12450:     logger.info("Pipeline completed successfully")
12451:     return 0
12452: else:
12453:     logger.warning("Pipeline completed with validation warnings")
12454:     return 0
12455:
12456: except ValidationError as e:
12457:     logger.error(f"Validation error: {e}")
12458:     return 1
12459: except ProcessingError as e:
12460:     logger.error(f"Processing error: {e}")
12461:     return 1
```

```
12462:     except SerializationError as e:
12463:         logger.error(f"Serialization error: {e}")
12464:         return 1
12465:     except Exception as e:
12466:         logger.error(f"Unexpected error in pipeline: {e}")
12467:         import traceback
12468:         traceback.print_exc()
12469:         return 1
12470:
12471: if __name__ == '__main__':
12472:     import argparse
12473:     import sys
12474:
12475:     # ConfiguraciÃ³n de argumentos CLI
12476:     parser = argparse.ArgumentParser(
12477:         description="Smart Policy Chunks Pipeline v3.0 - Industrial Grade",
12478:         formatter_class=argparse.RawDescriptionHelpFormatter,
12479:         epilog="""
12480: Examples:
12481:     python smart_policy_chunks_canonic_phase_one.py --input plan.pdf --output chunks.json
12482:     python smart_policy_chunks_canonic_phase_one.py --input plan.pdf --output chunks.json --save_full
12483:     """
12484:     )
12485:     parser.add_argument(
12486:         '--input',
12487:         type=str,
12488:         required=True,
12489:         help='Path to input policy document (required)'
12490:     )
12491:     parser.add_argument(
12492:         '--output',
12493:         type=str,
12494:         required=False,
12495:         default='output_chunks.json',
12496:         help='Path to output JSON file (default: output_chunks.json)'
12497:     )
12498:     parser.add_argument(
12499:         '--doc_id',
12500:         type=str,
12501:         required=False,
12502:         default='POL_PLAN_001',
12503:         help='Document identifier (default: POL_PLAN_001)'
12504:     )
12505:     parser.add_argument(
12506:         '--title',
12507:         type=str,
12508:         required=False,
12509:         default='Plan de Desarrollo',
12510:         help='Document title (default: Plan de Desarrollo)'
12511:     )
12512:     parser.add_argument(
12513:         '--max_chunks',
12514:         type=int,
12515:         required=False,
12516:         default=0,
12517:         help='Maximum number of chunks to generate (0 = unlimited, default: 0)'
```

```
12518:     )
12519:     parser.add_argument(
12520:         '--save_full',
12521:         action='store_true',
12522:         help='Save full version with complete text (creates *_full.json)'
12523:     )
12524:
12525:     # Parse arguments - allow argparse to exit normally on error
12526:     args = parser.parse_args()
12527:
12528:     # Execute main pipeline
12529:     exit_code = main(args)
12530:     sys.exit(exit_code)
12531:
12532:
12533:
12534: =====
12535: FILE: src/farfan_pipeline/processing/structural.py
12536: =====
12537:
12538: """
12539: Structural normalization with policy-awareness.
12540:
12541: Segments documents into policy-aware units.
12542: """
12543:
12544: from typing import Any
12545: from farfan_pipeline.core.calibration.decorators import calibrated_method
12546:
12547:
12548: class StructuralNormalizer:
12549:     """Policy-aware structural normalizer."""
12550:
12551:     @calibrated_method("farfan_core.processing.spc_ingestion.structural.StructuralNormalizer.normalize")
12552:     def normalize(self, raw_objects: dict[str, Any]) -> dict[str, Any]:
12553:         """
12554:             Normalize document structure with policy awareness.
12555:
12556:             Args:
12557:                 raw_objects: Raw parsed objects
12558:
12559:             Returns:
12560:                 Policy graph with structured sections
12561: """
12562:     policy_graph = {
12563:         "sections": [],
12564:         "policy_units": [],
12565:         "axes": [],
12566:         "programs": [],
12567:         "projects": [],
12568:         "years": [],
12569:         "territories": []
12570:     }
12571:
12572:     # Extract sections from pages
12573:     for page in raw_objects.get("pages", []):
```

```

12574:         text = page.get("text", "")
12575:
12576:         # Detect policy units
12577:         policy_units = self._detect_policy_units(text)
12578:         policy_graph["policy_units"].extend(policy_units)
12579:
12580:         # Create section
12581:         section = {
12582:             "text": text,
12583:             "page": page.get("page_num"),
12584:             "title": self._extract_title(text),
12585:             "area": None,
12586:             "eje": None,
12587:         }
12588:         policy_graph["sections"].append(section)
12589:
12590:         # Extract axes, programs, projects
12591:         for unit in policy_graph["policy_units"]:
12592:             if unit["type"] == "eje":
12593:                 policy_graph["axes"].append(unit["name"])
12594:             elif unit["type"] == "programa":
12595:                 policy_graph["programs"].append(unit["name"])
12596:             elif unit["type"] == "proyecto":
12597:                 policy_graph["projects"].append(unit["name"])
12598:
12599:         return policy_graph
12600:
12601: @calibrated_method("farfan_core.processing.spc_ingestion.structural.StructuralNormalizer._detect_policy_units")
12602: def _detect_policy_units(self, text: str) -> list[dict[str, Any]]:
12603:     """Detect policy units in text."""
12604:     units = []
12605:
12606:     # Simple keyword-based detection
12607:     keywords = {
12608:         "eje": ["eje", "pilar"],
12609:         "programa": ["programa"],
12610:         "proyecto": ["proyecto"],
12611:         "meta": ["meta"],
12612:         "indicador": ["indicador"],
12613:     }
12614:
12615:     for unit_type, keywords_list in keywords.items():
12616:         for keyword in keywords_list:
12617:             if keyword.lower() in text.lower():
12618:                 units.append({
12619:                     "type": unit_type,
12620:                     "name": f"{keyword.capitalize()} detected",
12621:                 })
12622:
12623:     return units
12624:
12625: @calibrated_method("farfan_core.processing.spc_ingestion.structural.StructuralNormalizer._extract_title")
12626: def _extract_title(self, text: str) -> str:
12627:     """Extract title from text."""
12628:     # Simple: first line or first N characters
12629:     lines = text.split("\n")

```

```
12630:         if lines:
12631:             return lines[0][:100]
12632:         return ""
12633:
12634:
12635:
12636: =====
12637: FILE: src/farfan_pipeline/processing/uncertainty_quantification.py
12638: =====
12639:
12640: """
12641: Bayesian Uncertainty Quantification for Aggregation
12642:
12643: This module implements rigorous uncertainty propagation for the aggregation pipeline.
12644: It provides confidence intervals, epistemic/aleatoric decomposition, and sensitivity analysis.
12645:
12646: Methods:
12647: - Bootstrap resampling for empirical confidence intervals
12648: - Bayesian error propagation for weighted averages
12649: - Variance decomposition (epistemic vs aleatoric uncertainty)
12650: - Monte Carlo sampling for non-linear aggregations
12651:
12652: References:
12653: - Efron, B. (1979). "Bootstrap methods: another look at the jackknife"
12654: - Saltelli, A. et al. (2020). "Five ways to ensure that models serve society"
12655: - Kendall, A. & Gal, Y. (2017). "What uncertainties do we need in Bayesian deep learning?"
12656: """
12657:
12658: from __future__ import annotations
12659:
12660: import logging
12661: from dataclasses import dataclass, field
12662: from typing import Any
12663:
12664: import numpy as np
12665: from scipy import stats
12666:
12667: logger = logging.getLogger(__name__)
12668:
12669:
12670: @dataclass(frozen=True)
12671: class UncertaintyMetrics:
12672:     """
12673:         Comprehensive uncertainty metrics for a score.
12674:
12675:     Attributes:
12676:         mean: Point estimate (expected value)
12677:         std: Standard deviation
12678:         variance: Variance
12679:         confidence_interval_95: 95% CI tuple (lower, upper)
12680:         confidence_interval_99: 99% CI tuple (lower, upper)
12681:         epistemic_uncertainty: Model uncertainty (reducible with more data)
12682:         aleatoric_uncertainty: Inherent data noise (irreducible)
12683:         coefficient_of_variation: std/mean (relative uncertainty)
12684:         skewness: Distribution asymmetry
12685:         kurtosis: Distribution tail heaviness
```

```
12686:         metadata: Additional diagnostic information
12687:         """
12688:         mean: float
12689:         std: float
12690:         variance: float
12691:         confidence_interval_95: tuple[float, float]
12692:         confidence_interval_99: tuple[float, float]
12693:         epistemic_uncertainty: float
12694:         aleatoric_uncertainty: float
12695:         coefficient_of_variation: float
12696:         skewness: float
12697:         kurtosis: float
12698:         metadata: dict[str, Any] = field(default_factory=dict)
12699:
12700:     def is_high_uncertainty(self, threshold: float = 0.3) -> bool:
12701:         """Check if coefficient of variation exceeds threshold."""
12702:         return self.coefficient_of_variation > threshold
12703:
12704:     def dominant_uncertainty_type(self) -> str:
12705:         """Identify whether epistemic or aleatoric dominates."""
12706:         if self.epistemic_uncertainty > self.aleatoric_uncertainty * 1.2:
12707:             return "epistemic"
12708:         elif self.aleatoric_uncertainty > self.epistemic_uncertainty * 1.2:
12709:             return "aleatoric"
12710:         else:
12711:             return "balanced"
12712:
12713:     def to_dict(self) -> dict[str, Any]:
12714:         """Convert to dictionary for serialization."""
12715:         return {
12716:             "mean": self.mean,
12717:             "std": self.std,
12718:             "variance": self.variance,
12719:             "confidence_interval_95": self.confidence_interval_95,
12720:             "confidence_interval_99": self.confidence_interval_99,
12721:             "epistemic_uncertainty": self.epistemic_uncertainty,
12722:             "aleatoric_uncertainty": self.aleatoric_uncertainty,
12723:             "coefficient_of_variation": self.coefficient_of_variation,
12724:             "skewness": self.skewness,
12725:             "kurtosis": self.kurtosis,
12726:             "is_high_uncertainty": self.is_high_uncertainty(),
12727:             "dominant_uncertainty": self.dominant_uncertainty_type(),
12728:             "metadata": self.metadata,
12729:         }
12730:
12731:
12732: class BootstrapAggregator:
12733:     """
12734:         Bootstrap resampling for uncertainty quantification in aggregation.
12735:
12736:         Bootstrap is a non-parametric method that works for any aggregation function,
12737:         including non-linear operations like Choquet integrals.
12738:     """
12739:
12740:     def __init__(self, n_samples: int = 1000, random_seed: int = 42):
12741:         """
```

```
12742:     Initialize bootstrap aggregator.
12743:
12744:     Args:
12745:         n_samples: Number of bootstrap resamples
12746:         random_seed: Fixed seed for reproducibility
12747:         """
12748:         self.n_samples = n_samples
12749:         self.rng = np.random.RandomState(random_seed)
12750:         logger.info(f"BootstrapAggregator initialized (n_samples={n_samples}, seed={random_seed})")
12751:
12752:     def bootstrap_weighted_average(
12753:         self,
12754:         scores: list[float],
12755:         weights: list[float] | None = None,
12756:     ) -> UncertaintyMetrics:
12757:         """
12758:             Compute weighted average with bootstrap confidence intervals.
12759:
12760:             Args:
12761:                 scores: List of input scores
12762:                 weights: Optional weights (default: uniform)
12763:
12764:             Returns:
12765:                 UncertaintyMetrics with full uncertainty decomposition
12766:
12767:             Raises:
12768:                 ValueError: If scores is empty or weights mismatch
12769:             """
12770:         if not scores:
12771:             raise ValueError("Cannot bootstrap empty score list")
12772:
12773:         scores_arr = np.array(scores, dtype=np.float64)
12774:         n = len(scores_arr)
12775:
12776:         if weights is None:
12777:             weights_arr = np.ones(n) / n
12778:         else:
12779:             if len(weights) != n:
12780:                 raise ValueError(f"Weight count {len(weights)} != score count {n}")
12781:             weights_arr = np.array(weights, dtype=np.float64)
12782:             weights_arr = weights_arr / np.sum(weights_arr) # Normalize
12783:
12784:         # Bootstrap resampling
12785:         resamples = np.zeros(self.n_samples)
12786:         for i in range(self.n_samples):
12787:             # Sample with replacement
12788:             indices = self.rng.choice(n, size=n, replace=True)
12789:             resample_scores = scores_arr[indices]
12790:             resample_weights = weights_arr[indices]
12791:             resample_weights = resample_weights / np.sum(resample_weights)
12792:
12793:             resamples[i] = np.sum(resample_scores * resample_weights)
12794:
12795:         # Compute statistics
12796:         mean = np.mean(resamples)
12797:         std = np.std(resamples, ddof=1)
```

```
12798:         variance = std ** 2
12799:
12800:     # Confidence intervals (percentile method)
12801:     ci_95_lower = np.percentile(resamples, 2.5)
12802:     ci_95_upper = np.percentile(resamples, 97.5)
12803:     ci_99_lower = np.percentile(resamples, 0.5)
12804:     ci_99_upper = np.percentile(resamples, 99.5)
12805:
12806:     # Uncertainty decomposition
12807:     # Epistemic: Uncertainty in the aggregation (reducible with more samples)
12808:     epistemic_uncertainty = std
12809:
12810:     # Aleatoric: Inherent variability in data (irreducible)
12811:     aleatoric_uncertainty = np.std(scores_arr, ddof=1)
12812:
12813:     # Coefficient of variation
12814:     cv = std / mean if mean != 0.0 else float('inf')
12815:
12816:     # Distribution shape
12817:     skewness = stats.skew(resamples)
12818:     kurtosis = stats.kurtosis(resamples)
12819:
12820:     logger.debug(
12821:         f"Bootstrap complete: mean={mean:.4f}, std={std:.4f}, "
12822:         f"CI95=[{ci_95_lower:.4f}, {ci_95_upper:.4f}]"
12823:     )
12824:
12825:     return UncertaintyMetrics(
12826:         mean=float(mean),
12827:         std=float(std),
12828:         variance=float(variance),
12829:         confidence_interval_95=(float(ci_95_lower), float(ci_95_upper)),
12830:         confidence_interval_99=(float(ci_99_lower), float(ci_99_upper)),
12831:         epistemic_uncertainty=float(epistemic_uncertainty),
12832:         aleatoric_uncertainty=float(aleatoric_uncertainty),
12833:         coefficient_of_variation=float(cv),
12834:         skewness=float(skewness),
12835:         kurtosis=float(kurtosis),
12836:         metadata={
12837:             "n_scores": n,
12838:             "n_bootstrap_samples": self.n_samples,
12839:             "original_scores_mean": float(np.mean(scores_arr)),
12840:             "original_scores_std": float(np.std(scores_arr, ddof=1)),
12841:         },
12842:     )
12843:
12844:
12845: class BayesianPropagation:
12846:     """
12847:         Bayesian error propagation for analytical uncertainty quantification.
12848:
12849:         For linear aggregations (weighted averages), we can compute uncertainty
12850:         analytically without Monte Carlo sampling.
12851:         """
12852:
12853:     @staticmethod
```

```

12854:     def propagate_weighted_average(
12855:         scores: list[float],
12856:         score_uncertainties: list[float],
12857:         weights: list[float] | None = None,
12858:     ) -> tuple[float, float]:
12859:         """
12860:             Analytical uncertainty propagation for weighted average.
12861:
12862:             Formula:
12863:                 If  $Y = \sum w_i * X_i$ , then
12864:                      $\text{Var}(Y) = \sum (w_i^2 * \text{Var}(X_i))$  [assuming independence]
12865:
12866:             Args:
12867:                 scores: List of score means
12868:                 score_uncertainties: List of score standard deviations
12869:                 weights: Optional weights (default: uniform)
12870:
12871:             Returns:
12872:                 Tuple of (mean, std) for aggregated score
12873:
12874:             Raises:
12875:                 ValueError: If input lists have mismatched lengths
12876:             """
12877:             n = len(scores)
12878:             if len(score_uncertainties) != n:
12879:                 raise ValueError(
12880:                     f"Score count {n} != uncertainty count {len(score_uncertainties)}"
12881:                 )
12882:
12883:             scores_arr = np.array(scores, dtype=np.float64)
12884:             uncertainties_arr = np.array(score_uncertainties, dtype=np.float64)
12885:
12886:             if weights is None:
12887:                 weights_arr = np.ones(n) / n
12888:             else:
12889:                 if len(weights) != n:
12890:                     raise ValueError(f"Weight count {len(weights)} != score count {n}")
12891:                 weights_arr = np.array(weights, dtype=np.float64)
12892:                 weights_arr = weights_arr / np.sum(weights_arr)
12893:
12894:             # Mean propagation
12895:             mean = np.sum(weights_arr * scores_arr)
12896:
12897:             # Variance propagation (assuming independence)
12898:             variance = np.sum((weights_arr ** 2) * (uncertainties_arr ** 2))
12899:             std = np.sqrt(variance)
12900:
12901:             logger.debug(
12902:                 f"Bayesian propagation: mean={mean:.4f}, std={std:.4f} "
12903:                 f"(from {n} inputs)"
12904:             )
12905:
12906:             return float(mean), float(std)
12907:
12908:
12909: class SensitivityAnalysis:

```

```
12910: """
12911:     Sensitivity analysis for identifying influential inputs.
12912:
12913:     Implements variance-based sensitivity analysis (Sobol indices) to quantify
12914:     the contribution of each input to output variance.
12915: """
12916:
12917: @staticmethod
12918: def compute_sobol_indices(
12919:     scores: list[float],
12920:     weights: list[float],
12921:     aggregation_func: callable = None,
12922:     n_samples: int = 1000,
12923:     random_seed: int = 42,
12924: ) -> dict[int, float]:
12925: """
12926:     Compute first-order Sobol indices for weighted average.
12927:
12928:     Sobol index S_i measures the fraction of output variance due to input i.
12929:
12930:     Args:
12931:         scores: Input scores
12932:         weights: Input weights
12933:         aggregation_func: Aggregation function (default: weighted average)
12934:         n_samples: Monte Carlo samples for estimation
12935:         random_seed: Reproducibility seed
12936:
12937:     Returns:
12938:         Dictionary mapping input index to Sobol index
12939:
12940:     Note:
12941:         For weighted average, Sobol index is proportional to weight^2 * var(score_i).
12942: """
12943: if aggregation_func is None:
12944:     # Default: weighted average
12945:     def aggregation_func(s, w):
12946:         return np.sum(np.array(s) * np.array(w))
12947:
12948:     n = len(scores)
12949:     scores_arr = np.array(scores, dtype=np.float64)
12950:     weights_arr = np.array(weights, dtype=np.float64)
12951:     weights_arr = weights_arr / np.sum(weights_arr)
12952:
12953:     rng = np.random.RandomState(random_seed)
12954:
12955:     # Baseline output variance
12956:     baseline_output = aggregation_func(scores_arr, weights_arr)
12957:
12958:     # Perturb each input and measure output variance
12959:     sobol_indices = {}
12960:     for i in range(n):
12961:         # Monte Carlo: sample score_i from distribution, hold others fixed
12962:         perturbed_outputs = []
12963:         for _ in range(n_samples):
12964:             perturbed_scores = scores_arr.copy()
12965:             # Perturb score_i (assume ±20% noise)
```

```
12966:             noise = rng.normal(0, 0.2 * abs(scores_arr[i]))
12967:             perturbed_scores[i] = scores_arr[i] + noise
12968:
12969:             perturbed_output = aggregation_func(perturbed_scores, weights_arr)
12970:             perturbed_outputs.append(perturbed_output)
12971:
12972:             # Sobol index = Var(E[Y|X_i]) / Var(Y)
12973:             # Approximated by: Var of perturbed outputs
12974:             variance_due_to_i = np.var(perturbed_outputs, ddof=1)
12975:             sobol_indices[i] = float(variance_due_to_i)
12976:
12977:             # Normalize so sum = 1.0
12978:             total_variance = sum(sobol_indices.values())
12979:             if total_variance > 0:
12980:                 sobol_indices = {k: v / total_variance for k, v in sobol_indices.items()}
12981:
12982:             logger.debug(f"Sobol indices: {sobol_indices}")
12983:             return sobol_indices
12984:
12985:
12986: def aggregate_with_uncertainty(
12987:     scores: list[float],
12988:     weights: list[float] | None = None,
12989:     n_bootstrap: int = 1000,
12990:     random_seed: int = 42,
12991: ) -> tuple[float, UncertaintyMetrics]:
12992:     """
12993:         Convenience function: Aggregate scores with full uncertainty quantification.
12994:
12995:     Args:
12996:         scores: Input scores
12997:         weights: Optional weights (default: uniform)
12998:         n_bootstrap: Number of bootstrap samples
12999:         random_seed: Reproducibility seed
13000:
13001:     Returns:
13002:         Tuple of (point_estimate, uncertainty_metrics)
13003:     """
13004:     bootstrapper = BootstrapAggregator(n_samples=n_bootstrap, random_seed=random_seed)
13005:     uncertainty = bootstrapper.bootstrap_weighted_average(scores, weights)
13006:
13007:     # Point estimate: use mean from bootstrap
13008:     point_estimate = uncertainty.mean
13009:
13010:     return point_estimate, uncertainty
13011:
13012:
13013: __all__ = [
13014:     "UncertaintyMetrics",
13015:     "BootstrapAggregator",
13016:     "BayesianPropagation",
13017:     "SensitivityAnalysis",
13018:     "aggregate_with_uncertainty",
13019: ]
```

```
13022:  
13023: =====  
13024: FILE: src/farfan_pipeline/question_context.py  
13025: =====  
13026:  
13027: from __future__ import annotations  
13028:  
13029: from collections.abc import Mapping as ABCMapping  
13030: from dataclasses import dataclass  
13031: from types import MappingProxyType  
13032: from typing import Any  
13033:  
13034:  
13035: def _freeze(obj: Any) -> Any: # noqa: ANN401 # type: ignore[misc]  
13036:     if isinstance(obj, ABCMapping):  
13037:         return MappingProxyType({k: _freeze(v) for k, v in obj.items()})  
13038:     if isinstance(obj, list | tuple):  
13039:         return tuple(_freeze(x) for x in obj)  
13040:     if isinstance(obj, set):  
13041:         return frozenset(_freeze(x) for x in obj)  
13042:     return obj  
13043:  
13044:  
13045: @dataclass(frozen=True, slots=True)  
13046: class QuestionContext:  
13047:     """Carries question requirements through entire pipeline (deep-immutable)."""  
13048:  
13049:     question_mapping: Any # type: ignore[misc]  
13050:     dnp_standards: ABCMapping[str, Any] # type: ignore[misc]  
13051:     required_evidence_types: tuple[str, ...]  
13052:     search_queries: tuple[str, ...]  
13053:     validation_criteria: ABCMapping[str, Any] # type: ignore[misc]  
13054:     traceability_id: str  
13055:  
13056:     def __post_init__(self) -> None: # type: ignore[misc]  
13057:         object.__setattr__(self, "question_mapping", _freeze(self.question_mapping))  
13058:         object.__setattr__(self, "dnp_standards", _freeze(self.dnp_standards))  
13059:         object.__setattr__(  
13060:             self, "required_evidence_types", tuple(self.required_evidence_types)  
13061:         )  
13062:         object.__setattr__(self, "search_queries", tuple(self.search_queries))  
13063:         object.__setattr__(  
13064:             self, "validation_criteria", _freeze(self.validation_criteria)  
13065:         )  
13066:  
13067:  
13068:  
13069: =====  
13070: FILE: src/farfan_pipeline/scoring/__init__.py  
13071: =====  
13072:  
13073: """Scoring module package.  
13074:  
13075: This package provides backward compatibility for code that imports  
13076: scoring classes from farfan_core.scoring.  
13077: """
```

```
13078:  
13079: from farfan_pipeline.analysis.scoring.scoring import (  
13080:     EvidenceStructureError,  
13081:     ModalityConfig,  
13082:     ModalityValidationException,  
13083:     QualityLevel,  
13084:     ScoredResult,  
13085:     ScoringError,  
13086:     ScoringModality,  
13087:     ScoringValidator,  
13088:     apply_rounding,  
13089:     apply_scoring,  
13090:     clamp,  
13091:     determine_quality_level,  
13092:     score_type_a,  
13093:     score_type_b,  
13094:     score_type_c,  
13095:     score_type_d,  
13096:     score_type_e,  
13097:     score_type_f,  
13098: )  
13099:  
13100: __all__ = [  
13101:     "EvidenceStructureError",  
13102:     "ModalityConfig",  
13103:     "ModalityValidationException",  
13104:     "QualityLevel",  
13105:     "ScoredResult",  
13106:     "ScoringError",  
13107:     "ScoringModality",  
13108:     "ScoringValidator",  
13109:     "apply_rounding",  
13110:     "apply_scoring",  
13111:     "clamp",  
13112:     "determine_quality_level",  
13113:     "score_type_a",  
13114:     "score_type_b",  
13115:     "score_type_c",  
13116:     "score_type_d",  
13117:     "score_type_e",  
13118:     "score_type_f",  
13119: ]  
13120:  
13121:  
13122:  
13123: ======  
13124: FILE: src/farfan_pipeline/scoring/scoring.py  
13125: ======  
13126:  
13127: """Scoring module - re-exports from analysis.scoring.scoring.  
13128:  
13129: This module provides backward compatibility for code that imports  
13130: scoring classes from farfan_core.scoring.scoring.  
13131: """  
13132:  
13133: from farfan_pipeline.analysis.scoring.scoring import (
```

```
13134:     EvidenceStructureError,
13135:     ModalityConfig,
13136:     ModalityValidationError,
13137:     QualityLevel,
13138:     ScoredResult,
13139:     ScoringError,
13140:     ScoringModality,
13141:     ScoringValidator,
13142:     apply_rounding,
13143:     apply_scoring,
13144:     clamp,
13145:     determine_quality_level,
13146:     score_type_a,
13147:     score_type_b,
13148:     score_type_c,
13149:     score_type_d,
13150:     score_type_e,
13151:     score_type_f,
13152: )
13153:
13154: __all__ = [
13155:     "EvidenceStructureError",
13156:     "ModalityConfig",
13157:     "ModalityValidationError",
13158:     "QualityLevel",
13159:     "ScoredResult",
13160:     "ScoringError",
13161:     "ScoringModality",
13162:     "ScoringValidator",
13163:     "apply_rounding",
13164:     "apply_scoring",
13165:     "clamp",
13166:     "determine_quality_level",
13167:     "score_type_a",
13168:     "score_type_b",
13169:     "score_type_c",
13170:     "score_type_d",
13171:     "score_type_e",
13172:     "score_type_f",
13173: ]
13174:
13175:
13176:
13177: =====
13178: FILE: src/farfan_pipeline/scripts/__init__.py
13179: =====
13180:
13181: """
13182: Runtime-facing console entrypoints for the SAAAAAA package.
13183:
13184: Modules in this package are designed to be executed via
13185: ``python -m farfan_core.scripts.<name>`` so that repository
13186: scripts no longer need to mutate ``sys.path`` at runtime.
13187: """
13188:
13189: from __future__ import annotations
```

```
13190:  
13191: __all__ = list[str] = []  
13192:  
13193:  
13194:  
13195: =====  
13196: FILE: src/farfan_pipeline/synchronization/__init__.py  
13197: =====  
13198:  
13199: from farfan_pipeline.synchronization.irrigation_synchronizer import ChunkMatrix  
13200:  
13201: __all__ = ["ChunkMatrix"]  
13202:  
13203:
```