

```
1: =====
2: F.A.R.F.A.N PIPELINE CODE AUDIT - BATCH 21
3: =====
4: Generated: 2025-12-07T06:17:28.008981
5: Files in this batch: 17
6: =====
7:
8:
9: =====
10: FILE: src/farfan_pipeline/synchronization/irrigation_synchronizer.py
11: =====
12:
13: """Irrigation synchronization and chunk matrix validation for policy analysis."""
14:
15: from farfan_pipeline.core.orchestrator.chunk_matrix_builder import (
16:     CHUNK_ID_PATTERN,
17:     DIMENSIONS,
18:     EXPECTED_CHUNK_COUNT,
19:     POLICY_AREAS,
20:     build_chunk_matrix,
21: )
22: from farfan_pipeline.core.types import ChunkData, PreprocessedDocument
23:
24:
25: class ChunkMatrix:
26:     """Validates and provides O(1) access to policy chunks organized by area \u2227 dimension.
27:
28:     Delegates validation to core.orchestrator.chunk_matrix_builder to ensure a
29:     single source of truth aligned with questionnaire_monolith.json.
30:     """
31:
32:     POLICY_AREAS = POLICY_AREAS
33:     DIMENSIONS = DIMENSIONS
34:     EXPECTED_CHUNK_COUNT = EXPECTED_CHUNK_COUNT
35:     CHUNK_ID_PATTERN = CHUNK_ID_PATTERN
36:
37:     def __init__(self, document: PreprocessedDocument) -> None:
38:         matrix, sorted_keys = build_chunk_matrix(document)
39:         self.chunks: dict[tuple[str, str], ChunkData] = matrix
40:         self._matrix = self.chunks # Backward compatibility for existing callers
41:         self.sorted_keys = tuple(sorted_keys)
42:
43:     def get_chunk(self, policy_area_id: str, dimension_id: str) -> ChunkData:
44:         """Retrieve chunk by policy area and dimension with O(1) lookup."""
45:         key = (policy_area_id, dimension_id)
46:         if key not in self.chunks:
47:             raise KeyError(f"Chunk not found for key: {policy_area_id}-{dimension_id}")
48:         return self.chunks[key]
49:
50:
51:
52: =====
53: FILE: src/farfan_pipeline/utils/__init__.py
54: =====
55:
56: """Utility functions and helpers for SAAAAAA system."""
```

```
57:
58:
59:
60: =====
61: FILE: src/farfan_pipeline/utils/contract_io.py
62: =====
63:
64: """
65: Contract I/O Envelope - Universal Metadata Wrapper
66: =====
67:
68: ContractEnvelope wraps every phase payload (input and output) with
69: universal metadata including schema_version, timestamp_utc, policy_unit_id,
70: content_digest, correlation_id, and deterministic event_id.
71:
72: This module provides the canonical wrapper for all phase I/O to ensure
73: consistent metadata handling across the pipeline.
74:
75: Author: Policy Analytics Research Unit
76: Version: 1.0.0
77: License: Proprietary
78: """
79:
80: from __future__ import annotations
81:
82: import hashlib
83: import json
84: from datetime import datetime, timezone
85: from typing import TYPE_CHECKING, Any
86:
87: from pydantic import BaseModel, Field, field_validator
88: from farfan_pipeline.core.calibration.decorators import calibrated_method
89:
90: if TYPE_CHECKING:
91:     from collections.abc import Mapping
92:
93: CANONICAL_SCHEMA_VERSION = "io-1.0"
94:
95:
96: def _canonical_json_bytes(obj: Any) -> bytes:
97:     """
98:         Convert object to canonical JSON bytes for deterministic hashing.
99:
100:     Args:
101:         obj: Object to serialize
102:
103:     Returns:
104:         Canonical JSON bytes (sorted keys, compact separators)
105:     """
106:     return json.dumps(
107:         obj,
108:         separators=(", ", ":"),  
sort_keys=True,  
ensure_ascii=False  
.encode("utf-8")
112:
```

```
113:
114: def sha256_hex(obj: Any) -> str:
115:     """
116:         Compute SHA-256 hex digest of object via canonical JSON.
117:
118:     Args:
119:         obj: Object to hash
120:
121:     Returns:
122:         64-character hexadecimal SHA-256 digest
123:
124:     Examples:
125:         >>> sha256_hex({"b": 2, "a": 1})
126:         'eed6d51ab37ca6df16a330c85094467efcab7b5746c0e02bc728a05069ede38b'
127:         >>> sha256_hex({"a": 1, "b": 2})  # Same despite key order
128:         'eed6d51ab37ca6df16a330c85094467efcab7b5746c0e02bc728a05069ede38b'
129:     """
130:     return hashlib.sha256(_canonical_json_bytes(obj)).hexdigest()
131:
132:
133: def utcnow_iso() -> str:
134:     """
135:         Get current UTC timestamp in ISO-8601 format with Z suffix.
136:
137:     Returns:
138:         ISO-8601 timestamp string with Z suffix
139:
140:     Examples:
141:         >>> ts = utcnow_iso()
142:         >>> ts.endswith('Z')
143:         True
144:         >>> 'T' in ts
145:         True
146:     """
147:     # Always Z-suffixed UTC; forbidden to use local time
148:     return datetime.now(timezone.utc).isoformat().replace("+00:00", "Z")
149:
150:
151: class ContractEnvelope(BaseModel):
152:     """
153:         Universal metadata wrapper for phase I/O payloads.
154:
155:         Wraps every phase input/output with consistent metadata for:
156:             - Schema versioning
157:             - Temporal tracking (UTC only)
158:             - Policy scope identification
159:             - Cryptographic verification
160:             - Request correlation
161:             - Deterministic event tracking
162:
163:         Fields:
164:             - schema_version: Logical envelope schema version
165:             - timestamp_utc: ISO-8601 Z-suffixed timestamp
166:             - policy_unit_id: Stable identifier selecting seeds & scope
167:             - correlation_id: Client-supplied request/run correlation
168:             - content_digest: SHA-256 over canonical JSON of 'payload'
```

```
169:     - event_id: Deterministic SHA-256 over (policy_unit_id, content_digest)
170:     - payload: The phase's typed deliverable/expectation
171:
172:     Examples:
173:         >>> payload = {"result": "success", "data": [1, 2, 3]}
174:         >>> env = ContractEnvelope.wrap(
175:             ...     payload,
176:             ...     policy_unit_id="PDM-001",
177:             ...     correlation_id="req-123"
178:             ... )
179:         >>> env.policy_unit_id
180:         'PDM-001'
181:         >>> len(env.content_digest)
182:         64
183:         >>> len(env.event_id)
184:         64
185:     """
186:
187:     schema_version: str = Field(default=CANONICAL_SCHEMA_VERSION)
188:     timestamp_utc: str = Field(default_factory=utcnow_iso)
189:     policy_unit_id: str = Field(min_length=1)
190:     correlation_id: str | None = Field(default=None)
191:     content_digest: str = Field(min_length=64, max_length=64)
192:     event_id: str = Field(min_length=64, max_length=64)
193:     payload: Mapping[str, Any] | Any
194:
195:     model_config = {"frozen": True, "extra": "forbid"}
196:
197:     @field_validator("timestamp_utc")
198:     @classmethod
199:     def _validate_utc(cls, v: str) -> str:
200:         """Validate timestamp is Z-suffixed UTC ISO-8601."""
201:         # Quick sanity: must end with Z and parse
202:         if not v.endswith("Z"):
203:             raise ValueError("timestamp_utc must be Z-suffixed UTC (ISO-8601)")
204:         # Let it raise if invalid
205:         datetime.fromisoformat(v.replace("Z", "+00:00"))
206:         return v
207:
208:     @classmethod
209:     def wrap(
210:         cls,
211:         payload: Any,
212:         *,
213:         policy_unit_id: str,
214:         correlation_id: str | None = None,
215:         schema_version: str = CANONICAL_SCHEMA_VERSION,
216:     ) -> ContractEnvelope:
217:         """
218:             Wrap a payload with universal metadata envelope.
219:
220:             Args:
221:                 payload: The phase deliverable/expectation to wrap
222:                 policy_unit_id: Policy unit identifier
223:                 correlation_id: Optional request correlation ID
224:                 schema_version: Schema version (default: io-1.0)
```

```
225:
226:     Returns:
227:         ContractEnvelope with computed digests and event ID
228:
229:     Examples:
230:         >>> payload = {"status": "ok"}
231:         >>> env = ContractEnvelope.wrap(payload, policy_unit_id="PDM-001")
232:         >>> env.payload["status"]
233:         'ok'
234:         >>> env.policy_unit_id
235:         'PDM-001'
236:         """
237:         digest = sha256_hex(payload)
238:         event_id = sha256_hex({"policy_unit_id": policy_unit_id, "digest": digest})
239:         return cls(
240:             schema_version=schema_version,
241:             policy_unit_id=policy_unit_id,
242:             correlation_id=correlation_id,
243:             content_digest=digest,
244:             event_id=event_id,
245:             payload=payload,
246:         )
247:
248:
249: if __name__ == "__main__":
250:     import doctest
251:
252:     # Run doctests
253:     print("Running doctests...")
254:     doctest.testmod(verbose=True)
255:
256:     # Minimal deterministic check
257:     print("\n" + "="*60)
258:     print("ContractEnvelope Integration Tests")
259:     print("=".*60)
260:
261:     print("\n1. Testing deterministic digest computation:")
262:     p = {"a": 1, "b": [2, 3]}
263:     e1 = ContractEnvelope.wrap(p, policy_unit_id="PU_123")
264:     e2 = ContractEnvelope.wrap({"b": [2, 3], "a": 1}, policy_unit_id="PU_123")
265:     assert e1.content_digest == e2.content_digest  # canonical JSON stable
266:     assert e1.event_id == e2.event_id
267:     print(f"\u234\u223 Digest is deterministic: {e1.content_digest[:16]}...")
268:     print(f"\u234\u223 Event ID is deterministic: {e1.event_id[:16]}...")
269:
270:     print("\n2. Testing envelope immutability:")
271:     try:
272:         e1.payload = {"modified": True}
273:         print("\u234\u227 FAILED: Envelope should be immutable")
274:     except Exception:
275:         print("\u234\u223 Envelope is immutable (frozen)")
276:
277:     print("\n3. Testing UTC timestamp validation:")
278:     assert e1.timestamp_utc.endswith('Z')
279:     print(f"\u234\u223 Timestamp is UTC: {e1.timestamp_utc}")
```

```
281:     print("\n4. Testing correlation ID:")
282:     e3 = ContractEnvelope.wrap(p, policy_unit_id="PU_123", correlation_id="corr-456")
283:     assert e3.correlation_id == "corr-456"
284:     print(f"    \u2192 Correlation ID: {e3.correlation_id}")
285:
286:     print("\n" + "="*60)
287:     print("ContractEnvelope doctest OK - All tests passed!")
288:     print("="*60)
289:
290:
291:
292: =====
293: FILE: src/farfan_pipeline/utils/contracts.py
294: =====
295:
296: """
297: CONTRACT DEFINITIONS - Frozen Data Shapes
298: =====
299:
300: TypedDict and Protocol definitions for API contracts across modules.
301: All data shapes must be versioned and adapters maintained for one release cycle.
302:
303: Purpose: Replace ad-hoc dicts with typed structures to prevent:
304: - unexpected keyword argument errors
305: - missing required positional arguments
306: - 'str' object has no attribute 'text' errors
307: - 'bool' object is not iterable
308: - unhashable type: 'dict' in sets
309:
310: Version 2.0 Enhancement:
311: - Pydantic-based contracts with strict validation (enhanced_contracts.py)
312: - Backward compatibility with V1 TypedDict contracts maintained
313: - Domain-specific exceptions and structured logging
314: - Cryptographic content verification and deterministic execution
315: """
316:
317: from __future__ import annotations
318:
319: from dataclasses import dataclass
320: from typing import (
321:     TYPE_CHECKING,
322:     Any,
323:     Literal,
324:     Protocol,
325:     TypedDict,
326: )
327:
328: if TYPE_CHECKING:
329:     from collections.abc import Iterable, Mapping, Sequence
330:     from pathlib import Path
331:
332: # =====
333: # V2 ENHANCED CONTRACTS - Pydantic-based with strict validation
334: # =====
335: # Import V2 contracts from enhanced_contracts module
336: # Use these for new code; V1 contracts maintained for backward compatibility
```

```
337: from farfan_pipeline.utils.enhanced_contracts import (
338:     # Pydantic Models
339:     AnalysisInputV2,
340:     AnalysisOutputV2,
341:     BaseContract,
342:     # Exceptions
343:     ContractValidationError,
344:     DataIntegrityError,
345:     DocumentMetadataV2,
346:     ExecutionContextV2,
347:     FlowCompatibilityError,
348:     ProcessedTextV2,
349:     # Utilities
350:     StructuredLogger,
351:     SystemConfigError,
352:     compute_content_digest,
353:     utc_now_iso,
354: )
355:
356: # =====
357: # DOCUMENT CONTRACTS - V1
358: # =====
359:
360: class DocumentMetadataV1(TypedDict, total=True):
361:     """Document metadata shape - all fields required."""
362:     file_path: str
363:     file_name: str
364:     num_pages: int
365:     file_size_bytes: int
366:     file_hash: str
367:
368: class DocumentMetadataV1Optional(TypedDict, total=False):
369:     """Optional document metadata fields."""
370:     pdf_metadata: dict[str, Any]
371:     author: str
372:     title: str
373:     creation_date: str
374:
375: class ProcessedTextV1(TypedDict, total=True):
376:     """Shape for processed text output."""
377:     raw_text: str
378:     normalized_text: str
379:     language: str
380:     encoding: str
381:
382: class ProcessedTextV1Optional(TypedDict, total=False):
383:     """Optional processed text fields."""
384:     sentences: list[str]
385:     sections: list[dict[str, Any]]
386:     tables: Mapping[str, Any]
387:
388: # =====
389: # ANALYSIS CONTRACTS - V1
390: # =====
391:
392: class AnalysisInputV1(TypedDict, total=True):
```

```
393:     """Required fields for analysis input - keyword-only."""
394:     text: str
395:     document_id: str
396:
397: class AnalysisInputV1Optional(TypedDict, total=False):
398:     """Optional fields for analysis input."""
399:     metadata: Mapping[str, Any]
400:     context: Mapping[str, Any]
401:     sentences: Sequence[str]
402:
403: class AnalysisOutputV1(TypedDict, total=True):
404:     """Shape for analysis output."""
405:     dimension: str
406:     category: str
407:     confidence: float
408:     matches: Sequence[str]
409:
410: class AnalysisOutputV1Optional(TypedDict, total=False):
411:     """Optional analysis output fields."""
412:     positions: Sequence[int]
413:     evidence: Sequence[str]
414:     warnings: Sequence[str]
415:
416: # =====
417: # EXECUTION CONTRACTS - V1
418: # =====
419:
420: class ExecutionContextV1(TypedDict, total=True):
421:     """Execution context for method invocation."""
422:     class_name: str
423:     method_name: str
424:     document_id: str
425:
426: class ExecutionContextV1Optional(TypedDict, total=False):
427:     """Optional execution context fields."""
428:     raw_text: str
429:     text: str
430:     metadata: Mapping[str, Any]
431:     tables: Mapping[str, Any]
432:     sentences: Sequence[str]
433:
434: # =====
435: # ERROR REPORTING CONTRACTS
436: # =====
437:
438: class ContractMismatchError(TypedDict, total=True):
439:     """Standard error shape for contract mismatches."""
440:     error_code: Literal["ERR_CONTRACT_MISMATCH"]
441:     stage: str
442:     function: str
443:     parameter: str
444:     expected_type: str
445:     got_type: str
446:     producer: str
447:     consumer: str
448:
```

```
449: # =====
450: # PROTOCOLS FOR PLUGGABLE BEHAVIOR
451: # =====
452:
453: class TextProcessorProtocol(Protocol):
454:     """Protocol for text processing components."""
455:
456:     @calibrated_method("farfan_core.utils.contracts.TextProcessorProtocol.normalize_unicode")
457:     def normalize_unicode(self, text: str) -> str:
458:         """Normalize unicode characters in text."""
459:         ...
460:
461:     @calibrated_method("farfan_core.utils.contracts.TextProcessorProtocol.segment_into_sentences")
462:     def segment_into_sentences(self, text: str) -> Sequence[str]:
463:         """Segment text into sentences."""
464:         ...
465:
466: class DocumentLoaderProtocol(Protocol):
467:     """Protocol for document loading components."""
468:
469:     @calibrated_method("farfan_core.utils.contracts.DocumentLoaderProtocol.load_pdf")
470:     def load_pdf(self, *, pdf_path: Path) -> DocumentMetadataV1:
471:         """Load PDF and return metadata - keyword-only params."""
472:         ...
473:
474:     @calibrated_method("farfan_core.utils.contracts.DocumentLoaderProtocol.validate_pdf")
475:     def validate_pdf(self, *, pdf_path: Path) -> bool:
476:         """Validate PDF file - keyword-only params."""
477:         ...
478:
479: class AnalyzerProtocol(Protocol):
480:     """Protocol for analysis components."""
481:
482:     def analyze(
483:         self,
484:         *,
485:         text: str,
486:         document_id: str,
487:         metadata: Mapping[str, Any] | None = None,
488:     ) -> AnalysisOutputV1:
489:         """Analyze text and return structured output - keyword-only params."""
490:         ...
491:
492: # =====
493: # VALUE OBJECTS (prevent .text on strings)
494: # =====
495:
496: @dataclass(frozen=True, slots=True)
497: class TextDocument:
498:     """Wrapper to prevent passing plain str where structured text is required."""
499:     text: str
500:     document_id: str
501:     metadata: Mapping[str, Any]
502:
503:     def __post_init__(self) -> None:
504:         """Validate that text is non-empty."""
```

```
505:         if not isinstance(self.text, str):
506:             raise TypeError(
507:                 f"ERR_CONTRACT_MISMATCH: text must be str, got {type(self.text).__name__}")
508:         )
509:     if not self.text:
510:         raise ValueError("ERR_CONTRACT_MISMATCH: text cannot be empty")
511:
512: @dataclass(frozen=True, slots=True)
513: class SentenceCollection:
514:     """Type-safe collection of sentences (prevents iteration bugs)."""
515:     sentences: tuple[str, ...] # Immutable and hashable
516:
517:     def __post_init__(self) -> None:
518:         """Validate sentences are strings."""
519:         if not all(isinstance(s, str) for s in self.sentences):
520:             raise TypeError(
521:                 "ERR_CONTRACT_MISMATCH: All sentences must be strings"
522:             )
523:
524:     def __iter__(self) -> Iterable[str]:
525:         """Make iterable."""
526:         return iter(self.sentences)
527:
528:     def __len__(self) -> int:
529:         """Return count."""
530:         return len(self.sentences)
531:
532: # =====
533: # SENTINEL VALUES (avoid None ambiguity)
534: # =====
535:
536: class _MissingSentinel:
537:     """Sentinel type for missing optional parameters."""
538:
539:     def __repr__(self) -> str:
540:         return "<MISSING>"
541:
542: MISSING: _MissingSentinel = _MissingSentinel()
543:
544: # =====
545: # RUNTIME VALIDATION HELPERS
546: # =====
547:
548: def validate_contract(
549:     value: Any,
550:     expected_type: type,
551:     *,
552:     parameter: str,
553:     producer: str,
554:     consumer: str,
555: ) -> None:
556:     """
557:     Validate value matches expected contract at runtime.
558:
559:     Raises TypeError with structured error message on mismatch.
560:     """
561:
```

```
561:     if not isinstance(value, expected_type):
562:         error_msg = (
563:             f"ERR_CONTRACT_MISMATCH["
564:             f"param='{parameter}', "
565:             f"expected={expected_type.__name__}, "
566:             f"got={type(value).__name__}, "
567:             f"producer={producer}, "
568:             f"consumer={consumer}"
569:             f"]"
570:         )
571:         raise TypeError(error_msg)
572:
573: def validate_mapping_keys(
574:     mapping: Mapping[str, Any],
575:     required_keys: Sequence[str],
576:     *,
577:     producer: str,
578:     consumer: str,
579: ) -> None:
580:     """
581:     Validate mapping contains required keys.
582:
583:     Raises KeyError with structured message on missing keys.
584:     """
585:     missing = [key for key in required_keys if key not in mapping]
586:     if missing:
587:         error_msg = (
588:             f"ERR_CONTRACT_MISMATCH["
589:             f"missing_keys={missing}, "
590:             f"producer={producer}, "
591:             f"consumer={consumer}"
592:             f"]"
593:         )
594:         raise KeyError(error_msg)
595:
596: def ensure_iterable_not_string(
597:     value: Any,
598:     *,
599:     parameter: str,
600:     producer: str,
601:     consumer: str,
602: ) -> None:
603:     """
604:     Validate value is iterable but NOT a string or bytes.
605:
606:     Prevents "'bool' object is not iterable" and "iterate string as tokens" bugs.
607:     """
608:     if isinstance(value, (str, bytes)):
609:         raise TypeError(
610:             f"ERR_CONTRACT_MISMATCH["
611:             f"param='{parameter}', "
612:             f"expected=Iterable (not str/bytes), "
613:             f"got={type(value).__name__}, "
614:             f"producer={producer}, "
615:             f"consumer={consumer}"
616:             f"]"
```

```
617:         )
618:
619:     try:
620:         iter(value)
621:     except TypeError as e:
622:         raise TypeError(
623:             f"ERR_CONTRACT_MISMATCH["
624:             f"param='{parameter}', "
625:             f"expected=Iterable, "
626:             f"got={type(value).__name__}, "
627:             f"producer={producer}, "
628:             f"consumer={consumer}"
629:             f"]"
630:         ) from e
631:
632: def ensure_hashable(
633:     value: Any,
634:     *,
635:     parameter: str,
636:     producer: str,
637:     consumer: str,
638: ) -> None:
639:     """
640:     Validate value is hashable (can be added to set or used as dict key).
641:
642:     Prevents "unhashable type: 'dict'" errors.
643:     """
644:     try:
645:         hash(value)
646:     except TypeError as e:
647:         raise TypeError(
648:             f"ERR_CONTRACT_MISMATCH["
649:             f"param='{parameter}', "
650:             f"expected=Hashable, "
651:             f"got={type(value).__name__} (unhashable), "
652:             f"producer={producer}, "
653:             f"consumer={consumer}"
654:             f"]"
655:         ) from e
656:
657:
658: # =====
659: # MODULE EXPORTS
660: # =====
661:
662: __all__ = [
663:     # V2 Enhanced Contracts (Pydantic-based) - RECOMMENDED FOR NEW CODE
664:     "AnalysisInputV2",
665:     "AnalysisOutputV2",
666:     "BaseContract",
667:     "DocumentMetadataV2",
668:     "ExecutionContextV2",
669:     "ProcessedTextV2",
670:     # V2 Exceptions
671:     "ContractValidationException",
672:     "DataIntegrityError",
```

```
673:     "FlowCompatibilityError",
674:     "SystemConfigError",
675:     # V2 Utilities
676:     "StructuredLogger",
677:     "compute_content_digest",
678:     "utc_now_iso",
679:     # V1 Contracts (TypedDict-based) - BACKWARD COMPATIBILITY
680:     "AnalysisInputV1",
681:     "AnalysisInputV1Optional",
682:     "AnalysisOutputV1",
683:     "AnalysisOutputV1Optional",
684:     "AnalyzerProtocol",
685:     "ContractMismatchError",
686:     "DocumentLoaderProtocol",
687:     "DocumentMetadataV1",
688:     "DocumentMetadataV1Optional",
689:     "ExecutionContextV1",
690:     "ExecutionContextV1Optional",
691:     "MISSING",
692:     "ProcessedTextV1",
693:     "ProcessedTextV1Optional",
694:     "SentenceCollection",
695:     "TextDocument",
696:     "TextProcessorProtocol",
697:     "ensure_hashable",
698:     "ensure_iterable_not_string",
699:     "validate_contract",
700:     "validate_mapping_keys",
701: ]
702:
703:
704:
705: =====
706: FILE: src/farfan_pipeline/utils/contracts_runtime.py
707: =====
708:
709: """
710: Runtime Contract Validation using Pydantic.
711:
712: This module provides runtime validators for all TypedDict contracts defined
713: in core_contracts.py. These validators enforce:
714: - Value bounds and constraints
715: - Required vs optional fields with strict validation
716: - Schema versioning for backward compatibility
717: - Round-trip serialization guarantees
718:
719: The validators mirror the TypedDict shapes exactly but add runtime enforcement.
720: Use these at public API boundaries and orchestrator edges.
721:
722: Version: 1.0.0
723: Schema Version Format: sem-{major}.{minor}
724: """
725:
726: from typing import Any
727:
728: from pydantic import BaseModel, ConfigDict, Field, field_validator
```

```
729: from farfan_pipeline.core.calibration.decorators import calibrated_method
730:
731: # =====
732: # CONFIGURATION
733: # =====
734:
735: class StrictModel(BaseModel):
736:     """Base model with strict configuration for all contract validators."""
737:
738:     model_config = ConfigDict(
739:         extra='forbid', # Refuse unknown fields
740:         validate_assignment=True,
741:         str_strip_whitespace=True,
742:         validate_default=True,
743:         populate_by_name=True, # Allow both field name and alias
744:     )
745:
746: # =====
747: # ANALYZER_ONE.PY CONTRACTS
748: # =====
749:
750: class SemanticAnalyzerInputModel(StrictModel):
751:     """Runtime validator for SemanticAnalyzerInputContract.
752:
753:     Validates:
754:     - text is non-empty
755:     - schema_version follows sem-X.Y pattern
756:     - segments is a list of strings
757:     - ontology_params is a valid dict
758:
759:     Example:
760:         >>> model = SemanticAnalyzerInputModel(
761:             ...      text="El plan de desarrollo municipal...",
762:             ...      schema_version="sem-1.0"
763:             ... )
764:
765:     text: str = Field(min_length=1, description="Document text to analyze")
766:     segments: list[str] = Field(
767:         default_factory=list,
768:         description="Pre-segmented text chunks"
769:     )
770:     ontology_params: dict[str, Any] = Field(
771:         default_factory=dict,
772:         description="Domain-specific ontology parameters"
773:     )
774:     schema_version: str = Field(
775:         default="sem-1.0",
776:         pattern=r"sem-\d+\.\d+$",
777:         description="Contract schema version"
778:     )
779:
780:     @field_validator('text')
781:     @classmethod
782:     def text_not_empty(cls, v: str) -> str:
783:         """Ensure text is not just whitespace."""
784:         if not v or not v.strip():
```

```
785:         raise ValueError("text must contain non-whitespace characters")
786:     return v
787:
788: class SemanticAnalyzerOutputModel(StrictModel):
789:     """Runtime validator for SemanticAnalyzerOutputContract."""
790:     semantic_cube: dict[str, Any] = Field(description="Semantic analysis results")
791:     coherence_score: float = Field(ge=0.0, le=1.0, description="Coherence metric")
792:     complexity_score: float = Field(ge=0.0, description="Complexity metric")
793:     domain_classification: dict[str, float] = Field(
794:         description="Domain probability distribution"
795:     )
796:     schema_version: str = Field(
797:         default="sem-1.0",
798:         pattern=r"^\d+\.\d+$"
799:     )
800:
801:     @field_validator('domain_classification')
802:     @classmethod
803:     def validate_probabilities(cls, v: dict[str, float]) -> dict[str, float]:
804:         """Ensure all domain probabilities are in [0, 1]."""
805:         for domain, prob in v.items():
806:             if not (0.0 <= prob <= 1.0):
807:                 raise ValueError(f"Probability for {domain} must be in [0, 1], got {prob}")
808:         return v
809:
810: # =====
811: # DERECK_BEACH.PY CONTRACTS
812: # =====
813:
814: class CDAFFrameworkInputModel(StrictModel):
815:     """Runtime validator for CDAFFrameworkInputContract."""
816:     document_text: str = Field(min_length=1, description="Document to analyze")
817:     plan_metadata: dict[str, Any] = Field(
818:         default_factory=dict,
819:         description="Metadata about the plan"
820:     )
821:     config: dict[str, Any] = Field(
822:         default_factory=dict,
823:         description="Framework configuration"
824:     )
825:     schema_version: str = Field(default="sem-1.0", pattern=r"^\d+\.\d+$")
826:
827: class CDAFFrameworkOutputModel(StrictModel):
828:     """Runtime validator for CDAFFrameworkOutputContract."""
829:     causal_mechanisms: list[dict[str, Any]] = Field(
830:         default_factory=list,
831:         description="Identified causal mechanisms"
832:     )
833:     evidential_tests: dict[str, Any] = Field(
834:         default_factory=dict,
835:         description="Statistical test results"
836:     )
837:     bayesian_inference: dict[str, Any] = Field(
838:         default_factory=dict,
839:         description="Bayesian analysis results"
840:     )
```

```
841:     audit_results: dict[str, Any] = Field(
842:         default_factory=dict,
843:         description="Audit findings and recommendations"
844:     )
845:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
846:
847: # =====
848: # FINANCIERO_VIABILIDAD_TABLAS.PY CONTRACTS
849: # =====
850:
851: class PDETAnalyzerInputModel(StrictModel):
852:     """Runtime validator for PDETAnalyzerInputContract."""
853:     document_content: str = Field(min_length=1, description="Document content")
854:     extract_tables: bool = Field(default=True, description="Whether to extract tables")
855:     config: dict[str, Any] = Field(default_factory=dict, description="Configuration")
856:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
857:
858: class PDETAnalyzerOutputModel(StrictModel):
859:     """Runtime validator for PDETAnalyzerOutputContract."""
860:     extracted_tables: list[dict[str, Any]] = Field(
861:         default_factory=list,
862:         description="Extracted financial tables"
863:     )
864:     financial_indicators: dict[str, float] = Field(
865:         default_factory=dict,
866:         description="Calculated financial metrics"
867:     )
868:     viability_score: float = Field(
869:         ge=0.0, le=1.0,
870:         description="Overall viability score"
871:     )
872:     quality_scores: dict[str, float] = Field(
873:         default_factory=dict,
874:         description="Quality assessment scores"
875:     )
876:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
877:
878: # =====
879: # TEORIA_CAMBIO.PY CONTRACTS
880: # =====
881:
882: class TeoriaCambioInputModel(StrictModel):
883:     """Runtime validator for TeoriaCambioInputContract."""
884:     document_text: str = Field(min_length=1, description="Document to analyze")
885:     strategic_goals: list[str] = Field(
886:         default_factory=list,
887:         description="Identified strategic goals"
888:     )
889:     config: dict[str, Any] = Field(default_factory=dict, description="Configuration")
890:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
891:
892: class TeoriaCambioOutputModel(StrictModel):
893:     """Runtime validator for TeoriaCambioOutputContract."""
894:     causal_dag: dict[str, Any] = Field(
895:         default_factory=dict,
896:         description="Causal directed acyclic graph"
```

```
897:     )
898:     validation_results: dict[str, Any] = Field(
899:         default_factory=dict,
900:         description="Model validation results"
901:     )
902:     monte_carlo_results: dict[str, Any] | None = Field(
903:         default=None,
904:         description="Monte Carlo simulation results"
905:     )
906:     graph_visualizations: list[dict[str, Any]] | None = Field(
907:         default=None,
908:         description="Graph visualization data"
909:     )
910:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
911:
912: # =====
913: # CONTRADICTION_DETECCION.PY CONTRACTS
914: # =====
915:
916: class ContradictionDetectorInputModel(StrictModel):
917:     """Runtime validator for ContradictionDetectorInputContract."""
918:     text: str = Field(min_length=1, description="Text to analyze for contradictions")
919:     plan_name: str = Field(min_length=1, description="Name of the plan")
920:     dimension: str | None = Field(
921:         default=None,
922:         description="PolicyDimension enum value"
923:     )
924:     config: dict[str, Any] = Field(default_factory=dict, description="Configuration")
925:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
926:
927: class ContradictionDetectorOutputModel(StrictModel):
928:     """Runtime validator for ContradictionDetectorOutputContract."""
929:     contradictions: list[dict[str, Any]] = Field(
930:         default_factory=list,
931:         description="Detected contradictions"
932:     )
933:     confidence_scores: dict[str, float] = Field(
934:         default_factory=dict,
935:         description="Confidence in each detection"
936:     )
937:     temporal_conflicts: list[dict[str, Any]] = Field(
938:         default_factory=list,
939:         description="Temporal inconsistencies"
940:     )
941:     severity_scores: dict[str, float] = Field(
942:         default_factory=dict,
943:         description="Severity ratings"
944:     )
945:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
946:
947: # =====
948: # EMBEDDING_POLICY.PY CONTRACTS
949: # =====
950:
951: class EmbeddingPolicyInputModel(StrictModel):
952:     """Runtime validator for EmbeddingPolicyInputContract."""
```

```
953:     text: str = Field(min_length=1, description="Text to embed")
954:     dimensions: list[str] = Field(
955:         default_factory=list,
956:         description="Policy dimensions to analyze"
957:     )
958:     embedding_model_config: dict[str, Any] = Field(
959:         default_factory=dict,
960:         description="Embedding model configuration",
961:         alias="model_config"
962:     )
963:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
964:
965: class EmbeddingPolicyOutputModel(StrictModel):
966:     """Runtime validator for EmbeddingPolicyOutputContract."""
967:     embeddings: list[list[float]] = Field(
968:         default_factory=list,
969:         description="Generated embeddings"
970:     )
971:     similarity_scores: dict[str, float] = Field(
972:         default_factory=dict,
973:         description="Similarity metrics"
974:     )
975:     bayesian_evaluation: dict[str, Any] = Field(
976:         default_factory=dict,
977:         description="Bayesian evaluation results"
978:     )
979:     policy_metrics: dict[str, float] = Field(
980:         default_factory=dict,
981:         description="Policy-specific metrics"
982:     )
983:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
984:
985: # =====
986: # SEMANTIC_CHUNKING_POLICY.PY CONTRACTS
987: # =====
988:
989: class SemanticChunkingInputModel(StrictModel):
990:     """Runtime validator for SemanticChunkingInputContract."""
991:     text: str = Field(min_length=1, description="Text to chunk")
992:     preserve_structure: bool = Field(
993:         default=True,
994:         description="Whether to preserve document structure"
995:     )
996:     config: dict[str, Any] = Field(default_factory=dict, description="Configuration")
997:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
998:
999: class SemanticChunkingOutputModel(StrictModel):
1000:     """Runtime validator for SemanticChunkingOutputContract."""
1001:     chunks: list[dict[str, Any]] = Field(
1002:         default_factory=list,
1003:         description="Semantic chunks"
1004:     )
1005:     causal_dimensions: dict[str, dict[str, Any]] = Field(
1006:         default_factory=dict,
1007:         description="Causal dimension analysis"
1008:     )
```

```
1009:     key_excerpts: dict[str, list[str]] = Field(
1010:         default_factory=dict,
1011:         description="Key excerpts by category"
1012:     )
1013:     summary: dict[str, Any] = Field(
1014:         default_factory=dict,
1015:         description="Summary statistics"
1016:     )
1017:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
1018:
1019: # =====
1020: # POLICY_PROCESSOR.PY CONTRACTS
1021: # =====
1022:
1023: class PolicyProcessorInputModel(StrictModel):
1024:     """Runtime validator for PolicyProcessorInputContract."""
1025:     data: Any = Field(description="Raw data to process")
1026:     text: str = Field(min_length=1, description="Text content")
1027:     sentences: list[str] = Field(
1028:         default_factory=list,
1029:         description="Pre-segmented sentences"
1030:     )
1031:     tables: list[dict[str, Any]] = Field(
1032:         default_factory=list,
1033:         description="Extracted tables"
1034:     )
1035:     config: dict[str, Any] = Field(default_factory=dict, description="Configuration")
1036:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
1037:
1038: class PolicyProcessorOutputModel(StrictModel):
1039:     """Runtime validator for PolicyProcessorOutputContract."""
1040:     processed_data: dict[str, Any] = Field(
1041:         default_factory=dict,
1042:         description="Processed results"
1043:     )
1044:     evidence_bundles: list[dict[str, Any]] = Field(
1045:         default_factory=list,
1046:         description="Evidence bundles"
1047:     )
1048:     bayesian_scores: dict[str, float] = Field(
1049:         default_factory=dict,
1050:         description="Bayesian scores"
1051:     )
1052:     matched_patterns: list[dict[str, Any]] = Field(
1053:         default_factory=list,
1054:         description="Matched patterns"
1055:     )
1056:     schema_version: str = Field(default="sem-1.0", pattern=r"^sem-\d+\.\d+$")
1057:
1058: # =====
1059: # EXPORTS
1060: # =====
1061:
1062: __all__ = [
1063:     # Base
1064:     'StrictModel',
```

```
1065:  
1066:     # Analyzer_one  
1067:     'SemanticAnalyzerInputModel',  
1068:     'SemanticAnalyzerOutputModel',  
1069:  
1070:     # derek_beach  
1071:     'CDAFFrameworkInputModel',  
1072:     'CDAFFrameworkOutputModel',  
1073:  
1074:     # financiero_viability_tablas  
1075:     'PDETAnalyzerInputModel',  
1076:     'PDETAnalyzerOutputModel',  
1077:  
1078:     # teoria_cambio  
1079:     'TeoriaCambioInputModel',  
1080:     'TeoriaCambioOutputModel',  
1081:  
1082:     # contradiction_deteccion  
1083:     'ContradictionDetectorInputModel',  
1084:     'ContradictionDetectorOutputModel',  
1085:  
1086:     # embedding_policy  
1087:     'EmbeddingPolicyInputModel',  
1088:     'EmbeddingPolicyOutputModel',  
1089:  
1090:     # semantic_chunking_policy  
1091:     'SemanticChunkingInputModel',  
1092:     'SemanticChunkingOutputModel',  
1093:  
1094:     # policy_processor  
1095:     'PolicyProcessorInputModel',  
1096:     'PolicyProcessorOutputModel',  
1097: ]  
1098:  
1099:  
1100:  
1101: =====  
1102: FILE: src/farfan_pipeline/utils/core_contracts.py  
1103: =====  
1104:  
1105: """  
1106: Core Module Contracts - Type-safe API boundaries for pure library modules.  
1107:  
1108: This module defines InputContract and OutputContract TypedDicts for each  
1109: core module to establish clear API boundaries and enable dependency injection.  
1110:  
1111: Architectural Principles:  
1112: - Core modules receive all data via InputContract parameters  
1113: - Core modules return data via OutputContract structures  
1114: - No I/O operations within core modules  
1115: - All I/O happens in orchestrator/factory.py  
1116: - Type-safe contracts with strict typing  
1117:  
1118: Version: 1.1.0  
1119: Schema Version: sem-1.0 (initial stable release)  
1120: Status: Active - Runtime validation available in contracts_runtime.py
```

```
1121: """
1122:
1123: from typing import Any, TypedDict
1124: from farfan_pipeline.core.calibration.decorators import calibrated_method
1125:
1126: try:
1127:     from typing import NotRequired # Python 3.11+
1128: except ImportError:
1129:     from typing_extensions import NotRequired # Python 3.9-3.10
1130:
1131: # =====
1132: # ANALYZER_ONE.PY CONTRACTS
1133: # =====
1134:
1135: class SemanticAnalyzerInputContract(TypedDict):
1136:     """Input contract for SemanticAnalyzer methods.
1137:
1138:     Example:
1139:         {
1140:             "text": "El plan de desarrollo municipal...",
1141:             "segments": ["Segment 1", "Segment 2"],
1142:             "ontology_params": {"domain": "municipal"}
1143:         }
1144:     """
1145:     text: str
1146:     segments: NotRequired[list[str]]
1147:     ontology_params: NotRequired[dict[str, Any]]
1148:
1149: class SemanticAnalyzerOutputContract(TypedDict):
1150:     """Output contract for SemanticAnalyzer methods."""
1151:     semantic_cube: dict[str, Any]
1152:     coherence_score: float
1153:     complexity_score: float
1154:     domain_classification: dict[str, float]
1155:
1156: # =====
1157: # DERECK_BEACH.PY CONTRACTS
1158: # =====
1159:
1160: class CDAFFrameworkInputContract(TypedDict):
1161:     """Input contract for CDAFFramework (Causal Deconstruction Audit Framework)."""
1162:     document_text: str
1163:     plan_metadata: dict[str, Any]
1164:     config: NotRequired[dict[str, Any]]
1165:
1166: class CDAFFrameworkOutputContract(TypedDict):
1167:     """Output contract for CDAFFramework."""
1168:     causal_mechanisms: list[dict[str, Any]]
1169:     evidential_tests: dict[str, Any]
1170:     bayesian_inference: dict[str, Any]
1171:     audit_results: dict[str, Any]
1172:
1173: # =====
1174: # FINANCIERO_VIABILIDAD_TABLAS.PY CONTRACTS
1175: # =====
1176:
```

```
1177: class PDETAnalyzerInputContract(TypedDict):
1178:     """Input contract for PDET (Programas de Desarrollo con Enfoque Territorial) Analyzer."""
1179:     document_content: str
1180:     extract_tables: NotRequired[bool]
1181:     config: NotRequired[dict[str, Any]]
1182:
1183: class PDETAnalyzerOutputContract(TypedDict):
1184:     """Output contract for PDET Analyzer."""
1185:     extracted_tables: list[dict[str, Any]]
1186:     financial_indicators: dict[str, float]
1187:     viability_score: float
1188:     quality_scores: dict[str, float]
1189:
1190: # =====
1191: # TEORIA_CAMBIO.PY CONTRACTS
1192: # =====
1193:
1194: class TeoriaCambioInputContract(TypedDict):
1195:     """Input contract for Theory of Change analysis."""
1196:     document_text: str
1197:     strategic_goals: NotRequired[list[str]]
1198:     config: NotRequired[dict[str, Any]]
1199:
1200: class TeoriaCambioOutputContract(TypedDict):
1201:     """Output contract for Theory of Change analysis."""
1202:     causal_dag: dict[str, Any]
1203:     validation_results: dict[str, Any]
1204:     monte_carlo_results: NotRequired[dict[str, Any]]
1205:     graph_visualizations: NotRequired[list[dict[str, Any]]]
1206:
1207: # =====
1208: # CONTRADICTION_DETECCION.PY CONTRACTS
1209: # =====
1210:
1211: class ContradictionDetectorInputContract(TypedDict):
1212:     """Input contract for PolicyContradictionDetector."""
1213:     text: str
1214:     plan_name: str
1215:     dimension: NotRequired[str] # PolicyDimension enum value
1216:     config: NotRequired[dict[str, Any]]
1217:
1218: class ContradictionDetectorOutputContract(TypedDict):
1219:     """Output contract for PolicyContradictionDetector."""
1220:     contradictions: list[dict[str, Any]]
1221:     confidence_scores: dict[str, float]
1222:     temporal_conflicts: list[dict[str, Any]]
1223:     severity_scores: dict[str, float]
1224:
1225: # =====
1226: # EMBEDDING_POLICY.PY CONTRACTS
1227: # =====
1228:
1229: class EmbeddingPolicyInputContract(TypedDict):
1230:     """Input contract for embedding-based policy analysis."""
1231:     text: str
1232:     dimensions: NotRequired[list[str]]
```

```
1233:     model_config: NotRequired[dict[str, Any]]
1234:
1235: class EmbeddingPolicyOutputContract(TypedDict):
1236:     """Output contract for embedding policy analysis."""
1237:     embeddings: list[list[float]]
1238:     similarity_scores: dict[str, float]
1239:     bayesian_evaluation: dict[str, Any]
1240:     policy_metrics: dict[str, float]
1241:
1242: # =====
1243: # SEMANTIC_CHUNKING_POLICY.PY CONTRACTS
1244: # =====
1245:
1246: class SemanticChunkingInputContract(TypedDict):
1247:     """Input contract for semantic chunking and policy document analysis."""
1248:     text: str
1249:     preserve_structure: NotRequired[bool]
1250:     config: NotRequired[dict[str, Any]]
1251:
1252: class SemanticChunkingOutputContract(TypedDict):
1253:     """Output contract for semantic chunking."""
1254:     chunks: list[dict[str, Any]]
1255:     causal_dimensions: dict[str, dict[str, Any]]
1256:     key_excerpts: dict[str, list[str]]
1257:     summary: dict[str, Any]
1258:
1259: # =====
1260: # POLICY_PROCESSOR.PY CONTRACTS
1261: # =====
1262:
1263: class PolicyProcessorInputContract(TypedDict):
1264:     """Input contract for IndustrialPolicyProcessor."""
1265:     data: Any
1266:     text: str
1267:     sentences: NotRequired[list[str]]
1268:     tables: NotRequired[list[dict[str, Any]]]
1269:     config: NotRequired[dict[str, Any]]
1270:
1271: class PolicyProcessorOutputContract(TypedDict):
1272:     """Output contract for IndustrialPolicyProcessor."""
1273:     processed_data: dict[str, Any]
1274:     evidence_bundles: list[dict[str, Any]]
1275:     bayesian_scores: dict[str, float]
1276:     matched_patterns: list[dict[str, Any]]
1277:
1278: # =====
1279: # SHARED DATA STRUCTURES
1280: # =====
1281:
1282: class DocumentData(TypedDict):
1283:     """Standard document data structure from farfan_core.core.orchestrator.
1284:
1285:     This is what the orchestrator/factory provides to core modules.
1286:     """
1287:     raw_text: str
1288:     sentences: list[str]
```

```
1289:     tables: list[dict[str, Any]]
1290:     metadata: dict[str, Any]
1291:
1292:     __all__ = [
1293:         # Analyzer_one
1294:         'SemanticAnalyzerInputContract',
1295:         'SemanticAnalyzerOutputContract',
1296:
1297:         # derek_beach
1298:         'CDAFFrameworkInputContract',
1299:         'CDAFFrameworkOutputContract',
1300:
1301:         # financiero_viability_tablas
1302:         'PDETAnalyzerInputContract',
1303:         'PDETAnalyzerOutputContract',
1304:
1305:         # teoria_cambio
1306:         'TeoriaCambioInputContract',
1307:         'TeoriaCambioOutputContract',
1308:
1309:         # contradiction_deteccion
1310:         'ContradictionDetectorInputContract',
1311:         'ContradictionDetectorOutputContract',
1312:
1313:         # embedding_policy
1314:         'EmbeddingPolicyInputContract',
1315:         'EmbeddingPolicyOutputContract',
1316:
1317:         # semantic_chunking_policy
1318:         'SemanticChunkingInputContract',
1319:         'SemanticChunkingOutputContract',
1320:
1321:         # policy_processor
1322:         'PolicyProcessorInputContract',
1323:         'PolicyProcessorOutputContract',
1324:
1325:         # Shared
1326:         'DocumentData',
1327:     ]
1328:
1329:
1330:
1331: =====
1332: FILE: src/farfan_pipeline/utils/coverage_gate.py
1333: =====
1334:
1335: #!/usr/bin/env python3
1336: """
1337: Coverage Enforcement Gate
1338: =====
1339: Enforces hard-fail at <555 methods threshold + audit.json emission
1340:
1341: Requirements:
1342: - Count all public methods across Producer classes
1343: - Generate audit.json with method counts and validation results
1344: - Hard-fail if total methods < 555
```

```
1345: - Include schema validation results
1346: """
1347:
1348: import ast
1349: import json
1350: import sys
1351: from datetime import datetime
1352: from pathlib import Path
1353: from farfan_pipeline.core.calibration.decorators import calibrated_method
1354:
1355:
1356: def count_methods_in_class(filepath: Path, class_name: str) -> tuple[list[str], dict[str, int]]:
1357:     """Count public and private methods in a class and return method names"""
1358:     if not filepath.exists():
1359:         return [], {"public": 0, "private": 0, "total": 0}
1360:
1361:     with open(filepath, encoding='utf-8') as f:
1362:         tree = ast.parse(f.read())
1363:
1364:     method_names = []
1365:     method_counts = {
1366:         "public": 0,
1367:         "private": 0,
1368:         "total": 0
1369:     }
1370:
1371:     for node in ast.walk(tree):
1372:         if isinstance(node, ast.ClassDef) and node.name == class_name:
1373:             for item in node.body:
1374:                 if isinstance(item, ast.FunctionDef):
1375:                     method_names.append(item.name)
1376:                     if not item.name.startswith('_'):
1377:                         method_counts["public"] += 1
1378:                     else:
1379:                         method_counts["private"] += 1
1380:                 method_counts["total"] += 1
1381:
1382:     return method_names, method_counts
1383:
1384: def validate_schema_exists(module_dir: Path) -> tuple[bool, list[str]]:
1385:     """Validate that schema files exist for a module"""
1386:     if not module_dir.exists():
1387:         return False, []
1388:
1389:     schema_files = list(module_dir.glob("*.schema.json"))
1390:     return len(schema_files) > 0, [f.name for f in schema_files]
1391:
1392: def count_file_methods(filepath: Path) -> tuple[int, int]:
1393:     """Count all public and total methods in a file"""
1394:     if not filepath.exists():
1395:         return 0, 0
1396:
1397:     with open(filepath, encoding='utf-8') as f:
1398:         try:
1399:             tree = ast.parse(f.read())
1400:             public_methods = 0
```

```
1401:         all_methods = 0
1402:
1403:         for node in ast.walk(tree):
1404:             if isinstance(node, ast.FunctionDef):
1405:                 all_methods += 1
1406:                 if not node.name.startswith('_'):
1407:                     public_methods += 1
1408:
1409:             return public_methods, all_methods
1410:     except Exception as e:
1411:         print(f"Error parsing {filepath}: {e}")
1412:     return 0, 0
1413:
1414: def count_all_methods() -> dict[str, any]:
1415:     """Count all methods across all modules and producers"""
1416:
1417:     # All files to analyze
1418:     files_to_analyze = [
1419:         "financiero_viability_tablas.py",
1420:         "Analyzer_one.py",
1421:         "contradiction_deteccion.py",
1422:         "embedding_policy.py",
1423:         "teoria_cambio.py",
1424:         "derek_beach.py",
1425:         "policy_processor.py",
1426:         "report_assembly.py",
1427:         "semantic_chunking_policy.py"
1428:     ]
1429:
1430:     # Producer classes to check
1431:     producers = {
1432:         "SemanticChunkingProducer": "semantic_chunking_policy.py",
1433:         "EmbeddingPolicyProducer": "embedding_policy.py",
1434:         "DerekBeachProducer": "derek_beach.py",
1435:         "ReportAssemblyProducer": "report_assembly.py"
1436:     }
1437:
1438:     results = {
1439:         "timestamp": datetime.now().isoformat(),
1440:         "files": {},
1441:         "producers": {},
1442:         "totals": {
1443:             "file_public_methods": 0,
1444:             "file_total_methods": 0,
1445:             "producer_methods": 0,
1446:             "threshold": 555,
1447:             "meets_threshold": False
1448:         },
1449:         "schema_validation": {},
1450:         "audit_status": "PENDING"
1451:     }
1452:
1453:     # Count file methods
1454:     print("=" * 80)
1455:     print("FILE METHOD COUNTS")
1456:     print("=" * 80)
```

```
1457:  
1458:     for filepath_str in files_to_analyze:  
1459:         filepath = Path(filepath_str)  
1460:         public_methods, total_methods = count_file_methods(filepath)  
1461:         results["files"][filepath_str] = {  
1462:             "public_methods": public_methods,  
1463:             "total_methods": total_methods  
1464:         }  
1465:         results["totals"]["file_public_methods"] += public_methods  
1466:         results["totals"]["file_total_methods"] += total_methods  
1467:         print(f"{filepath_str}:45} | {public_methods:4} public | {total_methods:4} total")  
1468:  
1469: # Count Producer methods  
1470: print("\n" + "=" * 80)  
1471: print("PRODUCER METHOD COUNTS")  
1472: print("=" * 80)  
1473:  
1474: for class_name, filepath in producers.items():  
1475:     methods, counts = count_methods_in_class(Path(filepath), class_name)  
1476:     results["producers"][class_name] = {  
1477:         "file": filepath,  
1478:         "methods": methods,  
1479:         "counts": counts,  
1480:         "public_methods": counts["public"]  
1481:     }  
1482:     results["totals"]["producer_methods"] += counts["public"]  
1483:     print(f"{class_name}:45} | {counts['public']:3} public | {counts['private']:3} private | {counts['total']:3} total")  
1484:  
1485: # Update meets_threshold  
1486: results["totals"]["meets_threshold"] = (  
1487:     results["totals"]["file_total_methods"] >= 555  
1488: )  
1489:  
1490: # Validate schemas  
1491: print("\n" + "=" * 80)  
1492: print("SCHEMA VALIDATION")  
1493: print("=" * 80)  
1494:  
1495: schema_modules = [  
1496:     "semantic_chunking_policy",  
1497:     "embedding_policy",  
1498:     "derek_beach",  
1499:     "report_assembly"  
1500: ]  
1501:  
1502: for module in schema_modules:  
1503:     module_dir = Path("schemas") / module  
1504:     has_schemas, schema_files = validate_schema_exists(module_dir)  
1505:     results["schema_validation"][module] = {  
1506:         "has_schemas": has_schemas,  
1507:         "schema_files": schema_files,  
1508:         "schema_count": len(schema_files)  
1509:     }  
1510:     status = "\u234\u223" if has_schemas else "\u234\u227"  
1511:     print(f"{module}:35} | {status} | {len(schema_files)} schemas")  
1512:
```

```
1513:     # Determine audit status
1514:     all_have_schemas = all(
1515:         v["has_schemas"] for v in results["schema_validation"].values()
1516:     )
1517:
1518:     if results["totals"]["meets_threshold"] and all_have_schemas:
1519:         results["audit_status"] = "PASS"
1520:     else:
1521:         results["audit_status"] = "FAIL"
1522:
1523:     return results
1524:
1525: def main() -> int:
1526:     """Main entry point"""
1527:     print("\n" + "=" * 80)
1528:     print("COVERAGE ENFORCEMENT GATE")
1529:     print("=" * 80 + "\n")
1530:
1531:     # Count all methods
1532:     results = count_all_methods()
1533:
1534:     # Print summary
1535:     print("\n" + "=" * 80)
1536:     print("SUMMARY")
1537:     print("=" * 80)
1538:     print(f"Total file methods: {results['totals']['file_total_methods'][:4]}")
1539:     print(f"Total public methods: {results['totals']['file_public_methods'][:4]}")
1540:     print(f"Producer methods: {results['totals']['producer_methods'][:4]}")
1541:     print(f"Threshold: {results['totals']['threshold'][:4]}")
1542:     print(f"Meets threshold: {results['totals']['meets_threshold']}")
1543:     print(f"All schemas present: {all(v['has_schemas'] for v in results['schema_validation'].values())}")
1544:     print(f"Audit status: {results['audit_status']}")
1545:
1546:     # Save audit.json
1547:     audit_path = Path("audit.json")
1548:     with open(audit_path, 'w', encoding='utf-8') as f:
1549:         json.dump(results, f, indent=2)
1550:
1551:     print(f"\n\234\223 Audit results saved to {audit_path}")
1552:
1553:     # Enforce hard-fail
1554:     if not results['totals']['meets_threshold']:
1555:         print("\n" + "=" * 80)
1556:         print("\u235\214 COVERAGE GATE FAILED")
1557:         print("=" * 80)
1558:         print(f"Required: {results['totals']['threshold']} methods")
1559:         print(f"Found: {results['totals']['file_total_methods']} methods")
1560:         print(f"Gap: {results['totals']['threshold'] - results['totals']['file_total_methods']} methods")
1561:         print("=" * 80 + "\n")
1562:
1563:     return 1
1564:
1565:     # Check schema validation
1566:     if not all(v['has_schemas'] for v in results['schema_validation'].values()):
1567:         print("\n" + "=" * 80)
1568:         print("\u235\214 SCHEMA VALIDATION FAILED")
1569:         print("=" * 80)
```

```
1569:         for module, validation in results['schema_validation'].items():
1570:             if not validation['has_schemas']:
1571:                 print(f"Missing schemas for: {module}")
1572:             print("=" * 80 + "\n")
1573:             return 1
1574:
1575:     print("\n" + "=" * 80)
1576:     print("â\234\223 COVERAGE GATE PASSED")
1577:     print("=" * 80)
1578:     print(f"All {results['totals']['file_total_methods']} methods accounted for")
1579:     print(f"{results['totals']['file_public_methods']} public methods available")
1580:     print(f"{results['totals']['producer_methods']} producer methods exposed")
1581:     print("All schema contracts validated")
1582:     print("=" * 80 + "\n")
1583:
1584:     return 0
1585:
1586: if __name__ == "__main__":
1587:     sys.exit(main())
1588:
1589:
1590:
1591: =====
1592: FILE: src/farfan_pipeline/utils/cpp_adapter.py
1593: =====
1594:
1595: """CPP to Orchestrator Adapter.
1596:
1597: This adapter converts Canon Policy Package (CPP) documents from the ingestion pipeline
1598: into the orchestrator's PreprocessedDocument format.
1599:
1600: Note: This is the canonical adapter implementation. SPC (Smart Policy Chunks) is the
1601: precursor to CPP.
1602:
1603: Design Principles:
1604: - Preserves complete provenance information
1605: - Orders chunks by text_span.start for deterministic ordering
1606: - Computes provenance_completeness metric
1607: - Provides prescriptive error messages on failure
1608: - Supports micro, meso, and macro chunk resolutions
1609: - Optional dependencies handled gracefully (pyarrow, structlog)
1610: """
1611:
1612: from __future__ import annotations
1613:
1614: import logging
1615: from datetime import datetime, timezone
1616: from types import MappingProxyType
1617: from typing import Any
1618:
1619: from farfan_pipeline.core.parameters import ParameterLoaderV2
1620: from farfan_pipeline.core.types import ChunkData, PreprocessedDocument, Provenance
1621:
1622: logger = logging.getLogger(__name__)
1623:
1624: _EMPTY_MAPPING = MappingProxyType({})
```

```
1625:  
1626:  
1627: class CPPAdapterError(Exception):  
1628:     """Raised when CPP to PreprocessedDocument conversion fails."""  
1629:  
1630:     pass  
1631:  
1632:  
1633: class CPPAdapter:  
1634:     """  
1635:         Adapter to convert CanonPolicyPackage (CPP output) to PreprocessedDocument.  
1636:  
1637:         This is the canonical adapter for the FARFAN pipeline, converting the rich  
1638:         CanonPolicyPackage data into the format expected by the orchestrator.  
1639:     """  
1640:  
1641:     def __init__(self, enable_runtime_validation: bool = True) -> None:  
1642:         """Initialize the CPP adapter.  
1643:  
1644:             Args:  
1645:                 enable_runtime_validation: Enable WiringValidator for runtime contract checking  
1646:             """  
1647:         self.logger = logging.getLogger(self.__class__.__name__)  
1648:  
1649:         # Initialize WiringValidator for runtime contract validation  
1650:         self.enable_runtime_validation = enable_runtime_validation  
1651:         if enable_runtime_validation:  
1652:             try:  
1653:                 from farfan_pipeline.core.wiring.validation import WiringValidator  
1654:  
1655:                 self.wiring_validator = WiringValidator()  
1656:                 self.logger.info(  
1657:                     "WiringValidator enabled for runtime contract checking"  
1658:                 )  
1659:             except ImportError:  
1660:                 self.logger.warning(  
1661:                     "WiringValidator not available. Runtime validation disabled."  
1662:                 )  
1663:             self.wiring_validator = None  
1664:         else:  
1665:             self.wiring_validator = None  
1666:  
1667:     def to_preprocessed_document(  
1668:         self, canon_package: Any, document_id: str  
1669:     ) -> PreprocessedDocument:  
1670:         """  
1671:             Convert CanonPolicyPackage to PreprocessedDocument.  
1672:  
1673:             Args:  
1674:                 canon_package: CanonPolicyPackage from ingestion  
1675:                 document_id: Unique document identifier  
1676:  
1677:             Returns:  
1678:                 PreprocessedDocument ready for orchestrator  
1679:  
1680:             Raises:
```

```
1681:             CPPAdapterError: If conversion fails or data is invalid
1682:
1683:     CanonPolicyPackage Expected Attributes:
1684:         Required:
1685:             - chunk_graph: ChunkGraph with .chunks dict
1686:             - chunk_graph.chunks: dict of chunk objects with .text and .text_span
1687:
1688:         Optional (handled with hasattr checks):
1689:             - schema_version: str (default: 'SPC-2025.1')
1690:             - quality_metrics: object with metrics like provenance_completeness,
1691:                 structural_consistency, boundary_f1, kpi_linkage_rate,
1692:                 budget_consistency_score, temporal_robustness, chunk_context_coverage
1693:             - policy_manifest: object with axes, programs, projects, years, territories
1694:             - metadata: dict with optional 'spc_rich_data' key
1695:
1696:         Chunk Optional Attributes (handled with hasattr checks):
1697:             - entities: list of entity objects with .text attribute
1698:             - time_facets: object with .years list
1699:             - budget: object with amount, currency, year, use, source attributes
1700:
1701:     """
1702:     self.logger.info(
1703:         f"Converting CanonPolicyPackage to PreprocessedDocument: {document_id}"
1704:     )
1705:
1706:     # === COMPREHENSIVE VALIDATION PHASE (H1.5) ===
1707:     # 6-layer validation for robust phase-one output processing
1708:
1709:     # V1: Validate canon_package exists
1710:     if not canon_package:
1711:         raise CPPAdapterError(
1712:             "canon_package is None or empty."
1713:             "Ensure ingestion completed successfully."
1714:         )
1715:
1716:     # V2: Validate document_id
1717:     if (
1718:         not document_id
1719:         or not isinstance(document_id, str)
1720:         or not document_id.strip()
1721:     ):
1722:         raise CPPAdapterError(
1723:             f"document_id must be a non-empty string."
1724:             f"Received: {repr(document_id)}"
1725:         )
1726:
1727:     # V3: Validate chunk_graph exists
1728:     if not hasattr(canon_package, "chunk_graph") or not canon_package.chunk_graph:
1729:         raise CPPAdapterError(
1730:             "canon_package must have a valid chunk_graph."
1731:             "Check that SmartChunkConverter produced valid output."
1732:         )
1733:
1734:     chunk_graph = canon_package.chunk_graph
1735:
1736:     # V4: Validate chunks dict is non-empty
1737:     if not chunk_graph.chunks:
```

```
1737:             raise CPPAdapterError(
1738:                 "chunk_graph.chunks is empty - no chunks to process. "
1739:                 "Minimum 1 chunk required from phase-one."
1740:             )
1741:
1742:             # V5: Validate individual chunks have required attributes
1743:             validation_failures = []
1744:             for chunk_id, chunk in chunk_graph.chunks.items():
1745:                 if not hasattr(chunk, "text"):
1746:                     validation_failures.append(
1747:                         f"Chunk {chunk_id}: missing 'text' attribute"
1748:                     )
1749:                 elif not chunk.text or not chunk.text.strip():
1750:                     validation_failures.append(
1751:                         f"Chunk {chunk_id}: text is empty or whitespace"
1752:                     )
1753:
1754:                 if not hasattr(chunk, "text_span"):
1755:                     validation_failures.append(
1756:                         f"Chunk {chunk_id}: missing 'text_span' attribute"
1757:                     )
1758:                 elif not hasattr(chunk.text_span, "start") or not hasattr(
1759:                     chunk.text_span, "end"
1760:                 ):
1761:                     validation_failures.append(
1762:                         f"Chunk {chunk_id}: invalid text_span (missing start/end)"
1763:                     )
1764:
1765:             # V6: Report validation failures with context
1766:             if validation_failures:
1767:                 failure_summary = "\n - ".join(validation_failures)
1768:                 raise CPPAdapterError(
1769:                     f"Chunk validation failed ({len(validation_failures)} errors):\n - {failure_summary}\n"
1770:                     f"Total chunks: {len(chunk_graph.chunks)}\n"
1771:                     f"This indicates SmartChunkConverter produced invalid output."
1772:                 )
1773:
1774:             # Sort chunks by document position for deterministic ordering
1775:             sorted_chunks = sorted(
1776:                 chunk_graph.chunks.values(),
1777:                 key=lambda c:
1778:                     c.text_span.start if hasattr(c, "text_span") and c.text_span else 0
1779:                 ),
1780:             )
1781:
1782:             # === PHASE 2 HARDENING: STRICT CARDINALITY & METADATA ===
1783:             # Enforce exactly 60 chunks for SPC/CPP canonical documents as per Jobfront 1
1784:             processing_mode = "chunked"
1785:             degradation_reason = None
1786:
1787:             if len(sorted_chunks) != 60:
1788:                 raise CPPAdapterError(
1789:                     f"Cardinality mismatch: Expected 60 chunks for 'chunked' processing mode, "
1790:                     f"but found {len(sorted_chunks)}. This is a critical violation of the "
1791:                     f"SPC canonical format."
1792:                 )
```

```
1793:  
1794:     # Enforce metadata integrity  
1795:     for idx, chunk in enumerate(sorted_chunks):  
1796:         if not hasattr(chunk, "policy_area_id") or not chunk.policy_area_id:  
1797:             raise CPPAdapterError(f"Missing policy_area_id in chunk {chunk.id}")  
1798:         if not hasattr(chunk, "dimension_id") or not chunk.dimension_id:  
1799:             raise CPPAdapterError(f"Missing dimension_id in chunk {chunk.id}")  
1800:         if not hasattr(chunk, "chunk_type") or not chunk.chunk_type:  
1801:             raise CPPAdapterError(f"Missing chunk_type in chunk {chunk.id}")  
1802:  
1803:     self.logger.info(f"Processing {len(sorted_chunks)} chunks")  
1804:  
1805:     # Build full text by concatenating chunks  
1806:     full_text_parts: list[str] = []  
1807:     sentences: list[dict[str, Any]] = []  
1808:     sentence_metadata: list[dict[str, Any]] = []  
1809:     tables: list[dict[str, Any]] = []  
1810:     chunk_index: dict[str, int] = {}  
1811:     chunk_summaries: list[dict[str, Any]] = []  
1812:     chunks_data: list[ChunkData] = []  
1813:  
1814:     # Track indices for building indexes  
1815:     term_index: dict[str, list[int]] = {}  
1816:     numeric_index: dict[str, list[int]] = {}  
1817:     temporal_index: dict[str, list[int]] = {}  
1818:     entity_index: dict[str, list[int]] = {}  
1819:  
1820:     # Track running offset that matches how full_text is built  
1821:     current_offset = 0  
1822:  
1823:     provenance_with_data = 0  
1824:  
1825:     for idx, chunk in enumerate(sorted_chunks):  
1826:         chunk_text = chunk.text  
1827:         chunk_start = current_offset  
1828:         chunk_index[chunk.id] = idx  
1829:  
1830:         # Add to full text  
1831:         full_text_parts.append(chunk_text)  
1832:  
1833:         # Create sentence entry (each chunk is represented as a sentence for orchestrator compatibility)  
1834:         sentences.append(  
1835:             {  
1836:                 "text": chunk_text,  
1837:                 "chunk_id": chunk.id,  
1838:                 "resolution": (  
1839:                     chunk.resolution.value.lower()  
1840:                     if hasattr(chunk, "resolution")  
1841:                     else None  
1842:                 ),  
1843:             },  
1844:         )  
1845:  
1846:         # Create chunk metadata for per-sentence tracking  
1847:         chunk_end = chunk_start + len(chunk_text)  
1848:
```

```
1849:     # CRITICAL: Preserve PA\227DIM metadata for Phase 2 question routing
1850:     extra_metadata = {
1851:         "chunk_id": chunk.id,
1852:         "policy_area_id": (
1853:             chunk.policy_area_id if hasattr(chunk, "policy_area_id") else None
1854:         ),
1855:         "dimension_id": (
1856:             chunk.dimension_id if hasattr(chunk, "dimension_id") else None
1857:         ),
1858:         "resolution": (
1859:             chunk.resolution.value.lower()
1860:             if hasattr(chunk, "resolution")
1861:             else None
1862:         ),
1863:     }
1864:
1865:     # Add facets if available
1866:     if hasattr(chunk, "policy_facets") and chunk.policy_facets:
1867:         extra_metadata["policy_facets"] = {
1868:             "axes": (
1869:                 chunk.policy_facets.axes
1870:                 if hasattr(chunk.policy_facets, "axes")
1871:                 else []
1872:             ),
1873:             "programs": (
1874:                 chunk.policy_facets.programs
1875:                 if hasattr(chunk.policy_facets, "programs")
1876:                 else []
1877:             ),
1878:             "projects": (
1879:                 chunk.policy_facets.projects
1880:                 if hasattr(chunk.policy_facets, "projects")
1881:                 else []
1882:             ),
1883:         }
1884:
1885:     if hasattr(chunk, "time_facets") and chunk.time_facets:
1886:         extra_metadata["time_facets"] = {
1887:             "years": (
1888:                 chunk.time_facets.years
1889:                 if hasattr(chunk.time_facets, "years")
1890:                 else []
1891:             ),
1892:             "periods": (
1893:                 chunk.time_facets.periods
1894:                 if hasattr(chunk.time_facets, "periods")
1895:                 else []
1896:             ),
1897:         }
1898:
1899:     if hasattr(chunk, "geo_facets") and chunk.geo_facets:
1900:         extra_metadata["geo_facets"] = {
1901:             "territories": (
1902:                 chunk.geo_facets.territories
1903:                 if hasattr(chunk.geo_facets, "territories")
1904:                 else []
```

```
1905:             ),
1906:             "regions": (
1907:                 chunk.geo_facets.regions
1908:                 if hasattr(chunk.geo_facets, "regions")
1909:                 else []
1910:             ),
1911:         }
1912:
1913:         sentence_metadata.append(
1914:             {
1915:                 "index": idx,
1916:                 "page_number": None,
1917:                 "start_char": chunk_start,
1918:                 "end_char": chunk_end,
1919:                 "extra": dict(extra_metadata),
1920:             }
1921:         )
1922:
1923:         chunk_summary = {
1924:             "id": chunk.id,
1925:             "resolution": (
1926:                 chunk.resolution.value.lower()
1927:                 if hasattr(chunk, "resolution")
1928:                 else None
1929:             ),
1930:             "text_span": {"start": chunk_start, "end": chunk_end},
1931:             "policy_area_id": extra_metadata["policy_area_id"],
1932:             "dimension_id": extra_metadata["dimension_id"],
1933:             "has_kpi": hasattr(chunk, "kpi") and chunk.kpi is not None,
1934:             "has_budget": hasattr(chunk, "budget") and chunk.budget is not None,
1935:             "confidence": {
1936:                 "layout": (
1937:                     getattr(
1938:                         chunk.confidence,
1939:                         "layout",
1940:                         ParameterLoaderV2.get(
1941:                             "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1942:                             "auto_param_L256_66",
1943:                             0.0,
1944:                         ),
1945:                     )
1946:                     if hasattr(chunk, "confidence")
1947:                     else ParameterLoaderV2.get(
1948:                         "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1949:                         "auto_param_L256_108",
1950:                         0.0,
1951:                     )
1952:                 ),
1953:                 "ocr": (
1954:                     getattr(
1955:                         chunk.confidence,
1956:                         "ocr",
1957:                         ParameterLoaderV2.get(
1958:                             "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1959:                             "auto_param_L257_60",
1960:                             0.0,
```

```

1961:                 ),
1962:                 )
1963:             if hasattr(chunk, "confidence")
1964:                 else ParameterLoaderV2.get(
1965:                     "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1966:                     "auto_param_L257_102",
1967:                     0.0,
1968:                 )
1969:             ),
1970:             "typing": (
1971:                 getattr(
1972:                     chunk.confidence,
1973:                     "typing",
1974:                     ParameterLoaderV2.get(
1975:                         "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1976:                         "auto_param_L258_66",
1977:                         0.0,
1978:                     ),
1979:                 )
1980:             if hasattr(chunk, "confidence")
1981:             else ParameterLoaderV2.get(
1982:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
1983:                 "auto_param_L258_108",
1984:                 0.0,
1985:             )
1986:         ),
1987:     },
1988: }
1989: chunk_summaries.append(chunk_summary)
1990:
1991: if hasattr(chunk, "provenance") and chunk.provenance:
1992:     provenance_with_data += 1
1993:     if (
1994:         not hasattr(chunk.provenance, "page_number")
1995:         or chunk.provenance.page_number is None
1996:     ):
1997:         raise CPPAdapterError(
1998:             f"Missing provenance.page_number in chunk {chunk.id}"
1999:         )
2000:     if (
2001:         not hasattr(chunk.provenance, "section_header")
2002:         or not chunk.provenance.section_header
2003:     ):
2004:         raise CPPAdapterError(
2005:             f"Missing provenance.section_header in chunk {chunk.id}"
2006:         )
2007: else:
2008:     raise CPPAdapterError(f"Missing provenance in chunk {chunk.id}")
2009:
2010: # Advance offset by chunk length + 1 space separator
2011: current_offset = chunk_end + 1
2012:
2013: # Extract entities for entity_index
2014: if hasattr(chunk, "entities") and chunk.entities:
2015:     for entity in chunk.entities:
2016:         entity_text = (

```

```

2017:             entity.text if hasattr(entity, "text") else str(entity)
2018:         )
2019:         if entity_text not in entity_index:
2020:             entity_index[entity_text] = []
2021:             entity_index[entity_text].append(idx)
2022:
2023:     # Extract temporal markers for temporal_index
2024:     if hasattr(chunk, "time_facets") and chunk.time_facets:
2025:         if hasattr(chunk.time_facets, "years") and chunk.time_facets.years:
2026:             for year in chunk.time_facets.years:
2027:                 year_key = str(year)
2028:                 if year_key not in temporal_index:
2029:                     temporal_index[year_key] = []
2030:                     temporal_index[year_key].append(idx)
2031:
2032:     # Extract budget for tables
2033:     if hasattr(chunk, "budget") and chunk.budget:
2034:         budget = chunk.budget
2035:         tables.append(
2036:             {
2037:                 "table_id": f"budget_{idx}",
2038:                 "label": f"Budget: {budget.source if hasattr(budget, 'source') else 'Unknown'}",
2039:                 "amount": getattr(budget, "amount", 0),
2040:                 "currency": getattr(budget, "currency", "COP"),
2041:                 "year": getattr(budget, "year", None),
2042:                 "use": getattr(budget, "use", None),
2043:                 "source": getattr(budget, "source", None),
2044:             }
2045:         )
2046:
2047:     # Create ChunkData object
2048:     chunk_type_value = chunk.chunk_type
2049:     if chunk_type_value not in [
2050:         "diagnostic",
2051:         "activity",
2052:         "indicator",
2053:         "resource",
2054:         "temporal",
2055:         "entity",
2056:     ]:
2057:         raise CPPAdapterError(
2058:             f"Invalid chunk_type '{chunk_type_value}' in chunk {chunk.id}"
2059:         )
2060:
2061:     chunks_data.append(
2062:         ChunkData(
2063:             id=idx,
2064:             text=chunk_text,
2065:             chunk_type=chunk_type_value,
2066:             sentences=[idx],
2067:             tables=(
2068:                 [len(tables) - 1]
2069:                 if hasattr(chunk, "budget") and chunk.budget
2070:                 else []
2071:             ),
2072:             start_pos=chunk_start,

```

```

2073:             end_pos=chunk_end,
2074:             confidence=(
2075:                 getattr(chunk.confidence, "overall", 1.0)
2076:                 if hasattr(chunk, "confidence")
2077:                 else 1.0
2078:             ),
2079:             edges_out=[], # Edges populated later if needed or from chunk_graph
2080:             policy_area_id=extra_metadata["policy_area_id"],
2081:             dimension_id=extra_metadata["dimension_id"],
2082:             provenance=(
2083:                 Provenance(
2084:                     page_number=chunk.provenance.page_number,
2085:                     section_header=getattr(
2086:                         chunk.provenance, "section_header", None
2087:                     ),
2088:                     bbox=getattr(chunk.provenance, "bbox", None),
2089:                     span_in_page=getattr(
2090:                         chunk.provenance, "span_in_page", None
2091:                     ),
2092:                     source_file=getattr(chunk.provenance, "source_file", None),
2093:                 )
2094:                 if hasattr(chunk, "provenance") and chunk.provenance
2095:                 else None
2096:             ),
2097:         ),
2098:     )
2099:
2100:     # Join full text
2101:     full_text = " ".join(full_text_parts)
2102:
2103:     if not full_text:
2104:         raise CPPAdapterError("Generated full_text is empty")
2105:
2106:     # Build document indexes
2107:     indexes = {
2108:         "term_index": {k: tuple(v) for k, v in term_index.items()},
2109:         "numeric_index": {k: tuple(v) for k, v in numeric_index.items()},
2110:         "temporal_index": {k: tuple(v) for k, v in temporal_index.items()},
2111:         "entity_index": {k: tuple(v) for k, v in entity_index.items()},
2112:     }
2113:
2114:     # Build metadata from canon_package
2115:     metadata_dict = {
2116:         "adapter_source": "CPPAdapter",
2117:         "schema_version": (
2118:             canon_package.schema_version
2119:             if hasattr(canon_package, "schema_version")
2120:             else "SPC-2025.1"
2121:         ),
2122:         "chunk_count": len(sorted_chunks),
2123:         "processing_mode": "chunked",
2124:         "chunks": chunk_summaries,
2125:     }
2126:
2127:     # Add quality metrics if available
2128:     if hasattr(canon_package, "quality_metrics") and canon_package.quality_metrics:

```

```
2129:     qm = canon_package.quality_metrics
2130:     metadata_dict["quality_metrics"] = {
2131:         "provenance_completeness": (
2132:             qm.provenance_completeness
2133:             if hasattr(qm, "provenance_completeness")
2134:             else ParameterLoaderV2.get(
2135:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2136:                 "auto_param_L328_117",
2137:                 0.0,
2138:             )
2139:         ),
2140:         "structural_consistency": (
2141:             qm.structural_consistency
2142:             if hasattr(qm, "structural_consistency")
2143:             else ParameterLoaderV2.get(
2144:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2145:                 "auto_param_L329_114",
2146:                 0.0,
2147:             )
2148:         ),
2149:         "boundary_f1": (
2150:             qm.boundary_f1
2151:             if hasattr(qm, "boundary_f1")
2152:             else ParameterLoaderV2.get(
2153:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2154:                 "auto_param_L330_81",
2155:                 0.0,
2156:             )
2157:         ),
2158:         "kpi_linkage_rate": (
2159:             qm.kpi_linkage_rate
2160:             if hasattr(qm, "kpi_linkage_rate")
2161:             else ParameterLoaderV2.get(
2162:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2163:                 "auto_param_L331_96",
2164:                 0.0,
2165:             )
2166:         ),
2167:         "budget_consistency_score": (
2168:             qm.budget_consistency_score
2169:             if hasattr(qm, "budget_consistency_score")
2170:             else ParameterLoaderV2.get(
2171:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2172:                 "auto_param_L332_120",
2173:                 0.0,
2174:             )
2175:         ),
2176:         "temporal_robustness": (
2177:             qm.temporal_robustness
2178:             if hasattr(qm, "temporal_robustness")
2179:             else ParameterLoaderV2.get(
2180:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2181:                 "auto_param_L333_105",
2182:                 0.0,
2183:             )
2184:         ),
```

```
2185:         "chunk_context_coverage": (
2186:             qm.chunk_context_coverage
2187:             if hasattr(qm, "chunk_context_coverage")
2188:             else ParameterLoaderV2.get(
2189:                 "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2190:                 "auto_param_L334_114",
2191:                 0.0,
2192:             )
2193:         ),
2194:     }
2195:
2196:     # Add policy manifest if available
2197:     if hasattr(canon_package, "policy_manifest") and canon_package.policy_manifest:
2198:         pm = canon_package.policy_manifest
2199:         metadata_dict["policy_manifest"] = {
2200:             "axes": pm.axes if hasattr(pm, "axes") else [],
2201:             "programs": pm.programs if hasattr(pm, "programs") else [],
2202:             "projects": pm.projects if hasattr(pm, "projects") else [],
2203:             "years": pm.years if hasattr(pm, "years") else [],
2204:             "territories": pm.territories if hasattr(pm, "territories") else [],
2205:         }
2206:
2207:     # Add SPC rich data if available in metadata
2208:     if hasattr(canon_package, "metadata") and canon_package.metadata:
2209:         if "spc_rich_data" in canon_package.metadata:
2210:             metadata_dict["spc_rich_data"] = canon_package.metadata["spc_rich_data"]
2211:
2212:         if len(sorted_chunks) > 0:
2213:             metadata_dict["provenance_completeness"] = provenance_with_data / len(
2214:                 sorted_chunks
2215:             )
2216:
2217:     metadata = MappingProxyType(metadata_dict)
2218:
2219:     # Detect language (default to Spanish for Colombian policy documents)
2220:     language = "es"
2221:
2222:     # Create PreprocessedDocument (canonical orchestrator dataclass)
2223:     preprocessed_doc = PreprocessedDocument(
2224:         document_id=document_id,
2225:         raw_text=full_text,
2226:         sentences=sentences,
2227:         tables=tables,
2228:         metadata=dict(metadata),
2229:         sentence_metadata=sentence_metadata,
2230:         indexes=indexes,
2231:         structured_text={
2232:             "full_text": full_text,
2233:             "sections": (),
2234:             "page_boundaries": (),
2235:         },
2236:         language=language,
2237:         ingested_at=datetime.now(timezone.utc),
2238:         full_text=full_text,
2239:         chunks=chunks_data,
2240:         chunk_index=chunk_index,
```

```

2241:         chunk_graph={
2242:             "chunks": {cid: chunk_index[cid] for cid in chunk_index},
2243:             "edges": list(getattr(chunk_graph, "edges", [])),
2244:         },
2245:         processing_mode=processing_mode,
2246:     )
2247:
2248:     self.logger.info(
2249:         f"Conversion complete: {len(sentences)} sentences, "
2250:         f"{len(tables)} tables, {len(entity_index)} entities indexed"
2251:     )
2252:
2253: # RUNTIME VALIDATION: Validate Adapter \206\222 Orchestrator contract
2254: if self.wiring_validator is not None:
2255:     self.logger.info("Validating Adapter \206\222 Orchestrator contract (runtime)")
2256:     try:
2257:         # Convert PreprocessedDocument to dict for validation
2258:         preprocessed_dict = {
2259:             "document_id": preprocessed_doc.document_id,
2260:             "sentence_metadata": preprocessed_doc.sentence_metadata,
2261:             "resolution_index": {}, # Placeholder, as it's not generated by the adapter
2262:             "provenance_completeness": metadata_dict.get(
2263:                 "provenance_completeness",
2264:                 ParameterLoaderV2.get(
2265:                     "farfan_core.utils.cpp_adapter.CPPAdapter.__init__",
2266:                     "auto_param_L397_92",
2267:                     0.0,
2268:                 ),
2269:             ),
2270:         }
2271:         self.wiring_validator.validate_adapter_to_orchestrator(
2272:             preprocessed_dict
2273:         )
2274:         self.logger.info("\234\223 Adapter \206\222 Orchestrator contract validation passed")
2275:     except Exception as e:
2276:         self.logger.error(
2277:             f"Adapter \206\222 Orchestrator contract validation failed: {e}"
2278:         )
2279:         raise ValueError(
2280:             f"Runtime contract violation at Adapter \206\222 Orchestrator boundary: {e}"
2281:         ) from e
2282:
2283:     return preprocessed_doc
2284:
2285:
2286: def adapt_cpp_to_orchestrator(
2287:     canon_package: Any, document_id: str
2288: ) -> PreprocessedDocument:
2289: """
2290: Convenience function to adapt CPP to PreprocessedDocument.
2291:
2292: Args:
2293:     canon_package: CanonPolicyPackage from ingestion
2294:     document_id: Unique document identifier
2295:
2296: Returns:

```

```
2297:     PreprocessedDocument for orchestrator
2298:
2299:     Raises:
2300:         CPPAdapterError: If conversion fails
2301:     """
2302:     adapter = CPPAdapter()
2303:     return adapter.to_preprocessed_document(canon_package, document_id)
2304:
2305:
2306: __all__ = [
2307:     "CPPAdapter",
2308:     "CPPAdapterError",
2309:     "adapt_cpp_to_orchestrator",
2310: ]
2311:
2312:
2313:
2314: =====
2315: FILE: src/farfan_pipeline/utils/determinism/__init__.py
2316: =====
2317:
2318: """Determinism utilities for reproducible runs."""
2319:
2320: from farfan_pipeline.utils.determinism.seeds import DeterministicContext, SeedFactory
2321:
2322: __all__ = [
2323:     "DeterministicContext",
2324:     "SeedFactory",
2325: ]
2326:
2327:
2328:
2329: =====
2330: FILE: src/farfan_pipeline/utils/determinism/seeds.py
2331: =====
2332:
2333: """Deterministic seed management for reproducible execution."""
2334:
2335: from __future__ import annotations
2336:
2337: import hashlib
2338: import os
2339: import random
2340: from dataclasses import dataclass
2341: from typing import TYPE_CHECKING
2342: from farfan_pipeline.core.calibration.decorators import calibrated_method
2343:
2344: if TYPE_CHECKING:
2345:     from collections.abc import Iterable
2346:
2347: try:
2348:     import numpy as np
2349:     NUMPY_AVAILABLE = True
2350: except ImportError: # pragma: no cover - optional dependency
2351:     np = None # type: ignore
2352:     NUMPY_AVAILABLE = False
```

```
2353:
2354: class SeedFactory:
2355:     """Factory that derives stable seeds from canonical metadata."""
2356:
2357:     DEFAULT_SALT = b"PDM_DETERMINISM_SALT_2025"
2358:
2359:     def __init__(self, salt: bytes | None = None) -> None:
2360:         self._salt = salt or self.DEFAULT_SALT
2361:
2362:     @calibrated_method("farfan_core.utils.determinism.seeds.SeedFactory.derive_seed")
2363:     def derive_seed(self, components: Iterable[str]) -> int:
2364:         """Derive a deterministic 32-bit seed from ordered components."""
2365:
2366:         material = "|".join(str(component) for component in components)
2367:         digest = hashlib.sha256(self._salt + material.encode("utf-8")).digest()
2368:         return int.from_bytes(digest[:4], byteorder="big")
2369:
2370:     @calibrated_method("farfan_core.utils.determinism.seeds.SeedFactory.derive_run_seed")
2371:     def derive_run_seed(self, questionnaire_hash: str, run_id: str) -> int:
2372:         """Derive run-wide seed based on questionnaire hash and run identifier."""
2373:
2374:         return self.derive_seed([questionnaire_hash, run_id])
2375:
2376:     @calibrated_method("farfan_core.utils.determinism.seeds.SeedFactory.configure_environment")
2377:     def configure_environment(self, seed: int) -> None:
2378:         """Configure deterministic state for Python and NumPy."""
2379:
2380:         os.environ["PYTHONHASHSEED"] = str(seed)
2381:         random.seed(seed)
2382:         if NUMPY_AVAILABLE and np is not None:
2383:             np.random.seed(seed)
2384:
2385:     @dataclass
2386:     class DeterministicContext:
2387:         """Deterministic execution context shared with all producers."""
2388:
2389:         questionnaire_hash: str
2390:         run_id: str
2391:         seed: int
2392:         numpy_rng: np.random.Generator | None = None
2393:
2394:         @calibrated_method("farfan_core.utils.determinism.seeds.DeterministicContext.apply")
2395:         def apply(self) -> None:
2396:             """Apply deterministic seeding across the runtime environment."""
2397:
2398:             os.environ["PYTHONHASHSEED"] = str(self.seed)
2399:             random.seed(self.seed)
2400:             if NUMPY_AVAILABLE and np is not None:
2401:                 self.numpy_rng = np.random.default_rng(self.seed)
2402:
2403:     @classmethod
2404:     def from_factory(
2405:         cls,
2406:         factory: SeedFactory,
2407:         questionnaire_hash: str,
2408:         run_id: str
```

```
2409:     ) -> DeterministicContext:
2410:         seed = factory.derive_run_seed(questionnaire_hash, run_id)
2411:         context = cls(questionnaire_hash=questionnaire_hash, run_id=run_id, seed=seed)
2412:         context.apply()
2413:         return context
2414:
2415:
2416: =====
2417: FILE: src/farfan_pipeline/utils/determinism_helpers.py
2418: =====
2419: =====
2420:
2421: """
2422: Determinism Helpers - Centralized Seeding and State Management
2423: =====
2424:
2425: Provides centralized determinism enforcement for the entire pipeline:
2426: - Stable seed derivation from policy_unit_id and correlation_id
2427: - Context manager for scoped deterministic execution
2428: - Controls random, numpy.random, and other stochastic libraries
2429:
2430: Author: Policy Analytics Research Unit
2431: Version: 1.0.0
2432: License: Proprietary
2433: """
2434:
2435: from __future__ import annotations
2436:
2437: import json
2438: import os
2439: import random
2440: from contextlib import contextmanager
2441: from dataclasses import dataclass
2442: from hashlib import sha256
2443: from typing import TYPE_CHECKING, Any
2444:
2445: import numpy as np
2446: from farfan_pipeline.core.calibration.decorators import calibrated_method
2447:
2448: if TYPE_CHECKING:
2449:     from collections.abc import Iterator
2450:
2451:
2452: def _seed_from(*parts: Any) -> int:
2453:     """
2454:     Derive a 32-bit seed from arbitrary parts via SHA-256.
2455:
2456:     Args:
2457:         *parts: Components to hash (will be JSON-serialized)
2458:
2459:     Returns:
2460:         32-bit integer seed suitable for random/numpy
2461:
2462:     Examples:
2463:         >>> s1 = _seed_from("PU_123", "corr-1")
2464:         >>> s2 = _seed_from("PU_123", "corr-1")
```

```
2465:         >>> s1 == s2
2466:         True
2467:         >>> s3 = _seed_from("PU_123", "corr-2")
2468:         >>> s1 != s3
2469:         True
2470:         """
2471:         raw = json.dumps(parts, sort_keys=True, separators=(",", ":"), ensure_ascii=False)
2472:         # 32-bit seed for numpy/py random
2473:         return int(sha256(raw.encode("utf-8")).hexdigest()[:8], 16)
2474:
2475:
2476: @dataclass(frozen=True)
2477: class Seeds:
2478:     """Container for seeds used in deterministic execution."""
2479:     py: int
2480:     np: int
2481:
2482:
2483: @contextmanager
2484: def deterministic(
2485:     policy_unit_id: str | None = None,
2486:     correlation_id: str | None = None
2487: ) -> Iterator[Seeds]:
2488:     """
2489:     Context manager for deterministic execution.
2490:
2491:     Sets seeds for Python's random and NumPy's random based on
2492:     policy_unit_id and correlation_id. Seeds are derived deterministically
2493:     via SHA-256 hashing.
2494:
2495:     Args:
2496:         policy_unit_id: Policy unit identifier (default: env var or "default")
2497:         correlation_id: Correlation identifier (default: env var or "run")
2498:
2499:     Yields:
2500:         Seeds object with py and np seed values
2501:
2502:     Examples:
2503:         >>> with deterministic("PU_123", "corr-1") as seeds:
2504:             ...     v1 = random.random()
2505:             ...     a1 = np.random.rand(3)
2506:             >>> with deterministic("PU_123", "corr-1") as seeds:
2507:                 ...     v2 = random.random()
2508:                 ...     a2 = np.random.rand(3)
2509:                 >>> v1 == v2  # Deterministic
2510:                 True
2511:                 >>> np.array_equal(a1, a2)  # Deterministic
2512:                 True
2513:                 """
2514:         base = policy_unit_id or os.getenv("POLICY_UNIT_ID", "default")
2515:         salt = correlation_id or os.getenv("CORRELATION_ID", "run")
2516:         s = _seed_from("fixed", base, salt)
2517:
2518:         # Set seeds for both random modules
2519:         random.seed(s)
2520:         np.random.seed(s)
```

```
2521:  
2522:     try:  
2523:         yield Seeds(py=s, np=s)  
2524:     finally:  
2525:         # Keep deterministic state; caller may reseed per-phase if needed  
2526:         pass  
2527:  
2528:  
2529: def create_deterministic_rng(seed: int) -> np.random.Generator:  
2530:     """  
2531:         Create a deterministic NumPy random number generator.  
2532:  
2533:         Use this for local RNG that doesn't affect global state.  
2534:  
2535:         Args:  
2536:             seed: Integer seed  
2537:  
2538:         Returns:  
2539:             NumPy Generator instance  
2540:  
2541:         Examples:  
2542:             >>> rng = create_deterministic_rng(42)  
2543:             >>> v1 = rng.random()  
2544:             >>> rng = create_deterministic_rng(42)  
2545:             >>> v2 = rng.random()  
2546:             >>> v1 == v2  
2547:             True  
2548:     """  
2549:     return np.random.default_rng(seed)  
2550:  
2551:  
2552: if __name__ == "__main__":  
2553:     import doctest  
2554:  
2555:     # Run doctests  
2556:     print("Running doctests...")  
2557:     doctest.testmod(verbose=True)  
2558:  
2559:     # Integration tests  
2560:     print("\n" + "="*60)  
2561:     print("Determinism Integration Tests")  
2562:     print("=".*60)  
2563:  
2564:     print("\n1. Testing seed derivation:")  
2565:     s1 = _seed_from("PU_123", "corr-1")  
2566:     s2 = _seed_from("PU_123", "corr-1")  
2567:     s3 = _seed_from("PU_123", "corr-2")  
2568:     assert s1 == s2  
2569:     assert s1 != s3  
2570:     print(f"  \u2193\234\223 Same inputs \u2193\206\222 same seed: {s1}")  
2571:     print(f"  \u2193\234\223 Different inputs \u2193\206\222 different seed: {s3}")  
2572:  
2573:     print("\n2. Testing deterministic context with random:")  
2574:     with deterministic("PU_123", "corr-1") as seeds1:  
2575:         a = random.random()  
2576:         b = random.randint(0, 100)
```

```
2577:     with deterministic("PU_123", "corr-1") as seeds2:
2578:         c = random.random()
2579:         d = random.randint(0, 100)
2580:         assert a == c
2581:         assert b == d
2582:         print(f"  \u2192 Python random is deterministic: {a:.6f}")
2583:         print(f"  \u2192 Python randint is deterministic: {b}")
2584:
2585:         print("\n3. Testing deterministic context with numpy:")
2586:         with deterministic("PU_123", "corr-1") as seeds:
2587:             arr1 = np.random.rand(3).tolist()
2588:             with deterministic("PU_123", "corr-1") as seeds:
2589:                 arr2 = np.random.rand(3).tolist()
2590:                 assert arr1 == arr2
2591:                 print(f"  \u2192 NumPy random is deterministic: {arr1}")
2592:
2593:         print("\n4. Testing local RNG generator:")
2594:         rng1 = create_deterministic_rng(42)
2595:         v1 = rng1.random()
2596:         rng2 = create_deterministic_rng(42)
2597:         v2 = rng2.random()
2598:         assert v1 == v2
2599:         print(f"  \u2192 Local RNG is deterministic: {v1:.6f}")
2600:
2601:     print("\n5. Testing different correlation IDs produce different results:")
2602:     with deterministic("PU_123", "corr-A"):
2603:         val_a = random.random()
2604:     with deterministic("PU_123", "corr-B"):
2605:         val_b = random.random()
2606:     assert val_a != val_b
2607:     print("  \u2192 Different correlation \u2192 different values")
2608:     print(f"    corr-A: {val_a:.6f}")
2609:     print(f"    corr-B: {val_b:.6f}")
2610:
2611:     print("\n" + "="*60)
2612:     print("Determinism doctest OK - All tests passed!")
2613:     print("="*60)
2614:
2615:
2616:
2617: =====
2618: FILE: src/farfan_pipeline/utils/deterministic_execution.py
2619: =====
2620:
2621: """
2622: Deterministic Execution Utilities - Production Grade
2623: =====
2624:
2625: Utilities for ensuring deterministic, reproducible execution across
2626: the policy analysis pipeline.
2627:
2628: Features:
2629: - Deterministic random seed management
2630: - UTC-only timestamp handling
2631: - Structured execution logging
2632: - Side-effect isolation
```

```
2633: - Reproducible event ID generation
2634:
2635: Author: Policy Analytics Research Unit
2636: Version: 1.0.0
2637: License: Proprietary
2638: """
2639:
2640: import hashlib
2641: import logging
2642: import random
2643: import time
2644: import uuid
2645: from collections.abc import Callable, Iterator
2646: from contextlib import contextmanager
2647: from datetime import datetime, timezone
2648: from typing import Any
2649:
2650: import numpy as np
2651:
2652: from farfan_pipeline.utils.enhanced_contracts import StructuredLogger, utc_now_iso
2653: from farfan_pipeline.core.calibration.decorators import calibrated_method
2654:
2655: # =====
2656: # DETERMINISTIC SEED MANAGEMENT
2657: # =====
2658:
2659: class DeterministicSeedManager:
2660:     """
2661:         Manages random seeds for deterministic execution.
2662:
2663:         All stochastic operations must use seeds managed by this class to ensure
2664:         reproducibility across runs.
2665:
2666:         Examples:
2667:             >>> manager = DeterministicSeedManager(base_seed=42)
2668:             >>> with manager.scoped_seed("operation1"):
2669:                 ...     value = random.random()
2670:             >>> # Seed is automatically restored after context
2671:     """
2672:
2673:     def __init__(self, base_seed: int = 42) -> None:
2674:         """
2675:             Initialize seed manager with base seed.
2676:
2677:             Args:
2678:                 base_seed: Master seed for all derived seeds
2679:         """
2680:         self.base_seed = base_seed
2681:         self._seed_counter = 0
2682:         self._initialize_seeds(base_seed)
2683:
2684:         @calibrated_method("farfan_core.utils.deterministic_execution.DeterministicSeedManager._initialize_seeds")
2685:         def _initialize_seeds(self, seed: int) -> None:
2686:             """Initialize all random number generators with deterministic seeds."""
2687:             random.seed(seed)
2688:             np.random.seed(seed)
```

```
2689:         # For reproducibility, also set hash seed
2690:         # Note: PYTHONHASHSEED should be set in environment for full determinism
2691:
2692:     @calibrated_method("farfan_core.utils.deterministic_execution.DeterministicSeedManager.get_derived_seed")
2693:     def get_derived_seed(self, operation_name: str) -> int:
2694:         """
2695:             Generate a deterministic seed for a specific operation.
2696:
2697:             Args:
2698:                 operation_name: Unique name for the operation
2699:
2700:             Returns:
2701:                 Deterministic integer seed derived from operation name and base seed
2702:
2703:             Examples:
2704:                 >>> manager = DeterministicSeedManager(42)
2705:                 >>> seed1 = manager.get_derived_seed("test")
2706:                 >>> seed2 = manager.get_derived_seed("test")
2707:                 >>> seed1 == seed2 # Deterministic
2708:                 True
2709:             """
2710:             # Use cryptographic hash for stable seed derivation
2711:             hash_input = f"{self.base_seed}:{operation_name}".encode()
2712:             hash_digest = hashlib.sha256(hash_input).digest()
2713:             # Convert first 4 bytes to int
2714:             return int.from_bytes(hash_digest[:4], byteorder='big')
2715:
2716:     @contextmanager
2717:     @calibrated_method("farfan_core.utils.deterministic_execution.DeterministicSeedManager.scoped_seed")
2718:     def scoped_seed(self, operation_name: str) -> Iterator[int]:
2719:         """
2720:             Context manager for scoped seed usage.
2721:
2722:             Sets seeds for the operation, then restores original state.
2723:
2724:             Args:
2725:                 operation_name: Unique name for the operation
2726:
2727:             Yields:
2728:                 Derived seed for this operation
2729:
2730:             Examples:
2731:                 >>> manager = DeterministicSeedManager(42)
2732:                 >>> with manager.scoped_seed("my_operation") as seed:
2733:                     ...      result = random.randint(0, 100)
2734:             """
2735:             # Save current state
2736:             random_state = random.getstate()
2737:             np_state = np.random.get_state()
2738:
2739:             # Set new seed
2740:             derived_seed = self.get_derived_seed(operation_name)
2741:             self._initialize_seeds(derived_seed)
2742:
2743:             try:
2744:                 yield derived_seed
```

```
2745:         finally:
2746:             # Restore state
2747:             random.setstate(random_state)
2748:             np.random.set_state(np_state)
2749:
2750:     @calibrated_method("farfan_core.utils.deterministic_execution.DeterministicSeedManager.get_event_id")
2751:     def get_event_id(self, operation_name: str, timestamp_utc: str | None = None) -> str:
2752:         """
2753:             Generate a reproducible event ID for an operation.
2754:
2755:             Args:
2756:                 operation_name: Operation name
2757:                 timestamp_utc: Optional UTC timestamp (ISO-8601); if None, uses current time
2758:
2759:             Returns:
2760:                 Deterministic event ID based on operation and timestamp
2761:
2762:             Examples:
2763:                 >>> manager = DeterministicSeedManager(42)
2764:                 >>> event_id = manager.get_event_id("test", "2024-01-01T00:00:00Z")
2765:                 >>> len(event_id)
2766:                 64
2767:             """
2768:             ts = timestamp_utc or utc_now_iso()
2769:             hash_input = f"{self.base_seed}:{operation_name}:{ts}".encode()
2770:             return hashlib.sha256(hash_input).hexdigest()
2771:
2772:
2773: # =====
2774: # DETERMINISTIC EXECUTION WRAPPER
2775: # =====
2776:
2777: class DeterministicExecutor:
2778:     """
2779:         Wraps functions to ensure deterministic execution with observability.
2780:
2781:         Features:
2782:             - Automatic seed management
2783:             - Structured logging of execution
2784:             - Latency tracking
2785:             - Error handling with event IDs
2786:
2787:         Examples:
2788:             >>> executor = DeterministicExecutor(base_seed=42, logger_name="test")
2789:             >>> @executor.deterministic(operation_name="my_func")
2790:             ... def my_function(x: int) -> int:
2791:                 ...     return x + random.randint(0, 10)
2792:             """
2793:
2794:     def __init__(
2795:         self,
2796:         base_seed: int = 42,
2797:         logger_name: str = "deterministic_executor",
2798:         enable_logging: bool = True
2799:     ) -> None:
2800:         """
```

```
2801:     Initialize deterministic executor.
2802:
2803:     Args:
2804:         base_seed: Master seed for all operations
2805:         logger_name: Logger name for structured logging
2806:         enable_logging: Whether to enable structured logging
2807:     """
2808:     self.seed_manager = DeterministicSeedManager(base_seed)
2809:     self.logger = StructuredLogger(logger_name) if enable_logging else None
2810:     self.enable_logging = enable_logging
2811:
2812:     def deterministic(
2813:         self,
2814:         operation_name: str,
2815:         log_inputs: bool = False,
2816:         log_outputs: bool = False
2817:     ) -> Callable:
2818:         """
2819:             Decorator to make a function deterministic with logging.
2820:
2821:             Args:
2822:                 operation_name: Unique name for this operation
2823:                 log_inputs: Whether to log input parameters
2824:                 log_outputs: Whether to log output values
2825:
2826:             Returns:
2827:                 Decorated function with deterministic execution
2828:         """
2829:     def decorator(func: Callable) -> Callable:
2830:         def wrapper(*args: Any, **kwargs: Any) -> Any:
2831:             # Generate correlation and event IDs
2832:             correlation_id = str(uuid.uuid4())
2833:             event_id = self.seed_manager.get_event_id(operation_name)
2834:
2835:             # Start timing
2836:             start_time = time.perf_counter()
2837:
2838:             # Execute with scoped seed
2839:             try:
2840:                 with self.seed_manager.scoped_seed(operation_name) as seed:
2841:                     result = func(*args, **kwargs)
2842:
2843:                     # Calculate latency
2844:                     latency_ms = (time.perf_counter() - start_time) * 1000
2845:
2846:                     # Log success
2847:                     if self.enable_logging and self.logger:
2848:                         log_data = {
2849:                             "event_id": event_id,
2850:                             "seed": seed,
2851:                             "latency_ms": latency_ms,
2852:                         }
2853:                         if log_inputs:
2854:                             log_data["inputs"] = str(args)[:100] # Truncate for safety
2855:                         if log_outputs:
2856:                             log_data["outputs"] = str(result)[:100]
```

```
2857:                     self.logger.log_execution(
2858:                         operation=operation_name,
2859:                         correlation_id=correlation_id,
2860:                         success=True,
2861:                         latency_ms=latency_ms,
2862:                         **log_data
2863:                     )
2864:
2865:
2866:             return result
2867:
2868:         except Exception as e:
2869:             # Calculate latency even on error
2870:             latency_ms = (time.perf_counter() - start_time) * 1000
2871:
2872:             # Log error
2873:             if self.enable_logging and self.logger:
2874:                 self.logger.log_execution(
2875:                     operation=operation_name,
2876:                     correlation_id=correlation_id,
2877:                     success=False,
2878:                     latency_ms=latency_ms,
2879:                     event_id=event_id,
2880:                     error=str(e)[:200] # Truncate for safety
2881:                 )
2882:
2883:             # Re-raise with event ID
2884:             raise RuntimeError(f"[{event_id}] {operation_name} failed: {e}") from e
2885:
2886:     return wrapper
2887:     return decorator
2888:
2889:
2890: # =====
2891: # UTC TIMESTAMP UTILITIES
2892: # =====
2893:
2894: def enforce_utc_now() -> datetime:
2895:     """
2896:     Get current UTC datetime.
2897:
2898:     Returns:
2899:         Current datetime in UTC timezone
2900:
2901:     Examples:
2902:         >>> dt = enforce_utc_now()
2903:         >>> dt.tzinfo is not None
2904:         True
2905:     """
2906:     return datetime.now(timezone.utc)
2907:
2908:
2909: def parse_utc_timestamp(timestamp_str: str) -> datetime:
2910:     """
2911:         Parse ISO-8601 timestamp and enforce UTC.
2912:
```

```
2913:     Args:
2914:         timestamp_str: ISO-8601 timestamp string
2915:
2916:     Returns:
2917:         Parsed datetime in UTC
2918:
2919:     Raises:
2920:         ValueError: If timestamp is not UTC or invalid format
2921:
2922:     Examples:
2923:         >>> dt = parse_utc_timestamp("2024-01-01T00:00:00Z")
2924:         >>> dt.year
2925:         2024
2926: """
2927: dt = datetime.fromisoformat(timestamp_str.replace('Z', '+00:00'))
2928:
2929: # Enforce UTC
2930: if dt.tzinfo is None or dt.utcoffset() != timezone.utc.utcoffset(None):
2931:     raise ValueError(f"Timestamp must be UTC: {timestamp_str}")
2932:
2933: return dt
2934:
2935:
2936: # =====
2937: # SIDE-EFFECT ISOLATION
2938: # =====
2939:
2940: @contextmanager
2941: def isolated_execution() -> Iterator[None]:
2942: """
2943:     Context manager to isolate side effects during execution.
2944:
2945:     Current isolation:
2946:         - Prevents print statements (captured and logged as warning)
2947:         - Future: file I/O restrictions, network restrictions
2948:
2949:     Yields:
2950:         None
2951:
2952:     Examples:
2953:         >>> with isolated_execution():
2954:             ...      # Code here has controlled side effects
2955:             ...      pass
2956: """
2957: # For now, minimal isolation - can be extended with more restrictions
2958: import io
2959: import sys
2960:
2961: # Capture stdout/stderr to detect violations
2962: old_stdout = sys.stdout
2963: old_stderr = sys.stderr
2964: stdout_capture = io.StringIO()
2965: stderr_capture = io.StringIO()
2966:
2967: try:
2968:     sys.stdout = stdout_capture
```

```
2969:         sys.stderr = stderr_capture
2970:         yield
2971:     finally:
2972:         sys.stdout = old_stdout
2973:         sys.stderr = old_stderr
2974:
2975:         # Log any captured output as warning (side effect violation)
2976:         if stdout_capture.getvalue():
2977:             logging.warning(
2978:                 "Side effect detected: stdout captured during isolated execution: %s",
2979:                 stdout_capture.getvalue()[:200]
2980:             )
2981:         if stderr_capture.getvalue():
2982:             logging.warning(
2983:                 "Side effect detected: stderr captured during isolated execution: %s",
2984:                 stderr_capture.getvalue()[:200]
2985:             )
2986:
2987:
2988: # =====
2989: # IN-SCRIPT TESTS
2990: # =====
2991:
2992: if __name__ == "__main__":
2993:     import doctest
2994:
2995:     # Run doctests
2996:     print("Running doctests...")
2997:     doctest.testmod(verbose=True)
2998:
2999:     # Additional tests
3000:     print("\n" + "="*60)
3001:     print("Deterministic Execution Tests")
3002:     print("=".*60)
3003:
3004:     # Test 1: Seed manager determinism
3005:     print("\n1. Testing seed manager determinism:")
3006:     manager1 = DeterministicSeedManager(42)
3007:     manager2 = DeterministicSeedManager(42)
3008:
3009:     seed1_a = manager1.get_derived_seed("test_op")
3010:     seed1_b = manager1.get_derived_seed("test_op")
3011:     seed2_a = manager2.get_derived_seed("test_op")
3012:
3013:     assert seed1_a == seed1_b == seed2_a, "Seeds must be deterministic"
3014:     print(f"\n2. Deterministic seeds: {seed1_a} == {seed1_b} == {seed2_a}")
3015:
3016:     # Test 2: Scoped seed restoration
3017:     print("\n2. Testing scoped seed restoration:")
3018:     manager = DeterministicSeedManager(42)
3019:
3020:     initial_value = random.random()
3021:     with manager.scoped_seed("temp_operation"):
3022:         _ = random.random() # Different value inside scope
3023:     restored_value = random.random()
3024:
```

```
3025:     # Reset and check if we can reproduce
3026:     manager._initialize_seeds(42)
3027:     reproduced_value = random.random()
3028:
3029:     print(f"    \u234\u223 Initial value: {initial_value:.6f}")
3030:     print(f"    \u234\u223 Reproduced value: {reproduced_value:.6f}")
3031:     assert abs(initial_value - reproduced_value) < 1e-10, "Seed restoration failed"
3032:     print("    \u234\u223 Seed restoration successful")
3033:
3034:     # Test 3: Deterministic executor
3035:     print("\n3. Testing deterministic executor:")
3036:     executor = DeterministicExecutor(base_seed=42, enable_logging=False)
3037:
3038:     @executor.deterministic(operation_name="test_function")
3039:     def sample_function(n: int) -> float:
3040:         return sum(random.random() for _ in range(n))
3041:
3042:     result1 = sample_function(5)
3043:
3044:     # Reset and run again
3045:     executor.seed_manager._initialize_seeds(42)
3046:     result2 = sample_function(5)
3047:
3048:     print(f"    \u234\u223 Result 1: {result1:.6f}")
3049:     print(f"    \u234\u223 Result 2: {result2:.6f}")
3050:     assert abs(result1 - result2) < 1e-10, "Deterministic execution failed"
3051:     print("    \u234\u223 Deterministic execution verified")
3052:
3053:     # Test 4: UTC enforcement
3054:     print("\n4. Testing UTC enforcement:")
3055:     utc_now = enforce_utc_now()
3056:     print(f"    \u234\u223 UTC now: {utc_now.isoformat()}")
3057:     assert utc_now.tzinfo is not None, "Must have timezone"
3058:
3059:     # Test 5: Event ID reproducibility
3060:     print("\n5. Testing event ID reproducibility:")
3061:     manager = DeterministicSeedManager(42)
3062:     event_id1 = manager.get_event_id("operation", "2024-01-01T00:00:00Z")
3063:     event_id2 = manager.get_event_id("operation", "2024-01-01T00:00:00Z")
3064:     assert event_id1 == event_id2, "Event IDs must be reproducible"
3065:     print(f"    \u234\u223 Event ID: {event_id1[:16]}...")
3066:     print("    \u234\u223 Event ID reproducibility verified")
3067:
3068:     print("\n" + "="*60)
3069:     print("All tests passed!")
3070:     print("="*60)
3071:
3072:
3073:
3074: =====
3075: FILE: src/farfan_pipeline/utils/domain_errors.py
3076: =====
3077:
3078: from farfan_pipeline.core.calibration.decorators import calibrated_method
3079: """
3080: Domain-Specific Exceptions - Contract Violation Errors
```

```
3081: =====
3082:
3083: Provides domain-specific exception hierarchy for contract violations.
3084:
3085: Exception Hierarchy:
3086:     ContractViolationError (base)
3087:         \224\234\224\200\224\200 DataContractError (data/payload violations)
3088:             \224\224\224\200\224\200 SystemContractError (system/configuration violations)
3089:
3090: Author: Policy Analytics Research Unit
3091: Version: 1.0.0
3092: License: Proprietary
3093: """
3094:
3095:
3096: class ContractViolationError(Exception):
3097:     """
3098:         Base exception for all contract violations.
3099:
3100:     Use this as the base class for specific contract violation types.
3101:
3102:     Examples:
3103:         >>> try:
3104:             ...      raise ContractViolationError("Contract violated")
3105:             ... except ContractViolationError as e:
3106:                 ...      print(f"Caught: {e}")
3107:             Caught: Contract violated
3108: """
3109: pass
3110:
3111:
3112: class DataContractError(ContractViolationError):
3113:     """
3114:         Exception for data/payload contract violations.
3115:
3116:         Raised when:
3117:             - Payload schema is invalid
3118:             - Required fields are missing
3119:             - Field values are out of range
3120:             - Data integrity checks fail (e.g., digest mismatch)
3121:
3122:         Examples:
3123:             >>> try:
3124:                 ...      raise DataContractError("Invalid payload schema")
3125:                 ... except DataContractError as e:
3126:                     ...      print(f"Data error: {e}")
3127:             Data error: Invalid payload schema
3128: """
3129: pass
3130:
3131:
3132: class SystemContractError(ContractViolationError):
3133:     """
3134:         Exception for system/configuration contract violations.
3135:
3136:         Raised when:
```

```
3137:     - System configuration is invalid
3138:     - Required resources are unavailable
3139:     - Environment preconditions are not met
3140:     - Infrastructure failures occur
3141:
3142: Examples:
3143:     >>> try:
3144:         ...      raise SystemContractError("Configuration missing")
3145:         ... except SystemContractError as e:
3146:             ...      print(f"System error: {e}")
3147:             System error: Configuration missing
3148: """
3149: pass
3150:
3151:
3152: if __name__ == "__main__":
3153:     import doctest
3154:
3155:     # Run doctests
3156:     print("Running doctests...")
3157:     doctest.testmod(verbose=True)
3158:
3159:     # Integration tests
3160:     print("\n" + "="*60)
3161:     print("Domain Exceptions Integration Tests")
3162:     print("="*60)
3163:
3164:     print("\n1. Testing exception hierarchy:")
3165:     assert issubclass(DataContractError, ContractViolationError)
3166:     assert issubclass(SystemContractError, ContractViolationError)
3167:     print("    \u234\223 DataContractError inherits from ContractViolationError")
3168:     print("    \u234\223 SystemContractError inherits from ContractViolationError")
3169:
3170:     print("\n2. Testing exception catching:")
3171:     try:
3172:         raise DataContractError("Test data error")
3173:     except ContractViolationError as e:
3174:         assert isinstance(e, DataContractError)
3175:         print("    \u234\223 DataContractError caught as ContractViolationError")
3176:
3177:     try:
3178:         raise SystemContractError("Test system error")
3179:     except ContractViolationError as e:
3180:         assert isinstance(e, SystemContractError)
3181:         print("    \u234\223 SystemContractError caught as ContractViolationError")
3182:
3183:     print("\n3. Testing specific exception catching:")
3184:     try:
3185:         raise DataContractError("Payload validation failed")
3186:     except DataContractError as e:
3187:         assert str(e) == "Payload validation failed"
3188:         print("    \u234\223 DataContractError caught specifically")
3189:
3190:     try:
3191:         raise SystemContractError("Config file not found")
3192:     except SystemContractError as e:
```

```
3193:         assert str(e) == "Config file not found"
3194:         print("    \u2193\233 SystemContractError caught specifically")
3195:
3196:     print("\n4. Testing error differentiation:")
3197:     errors = []
3198:
3199:     try:
3200:         raise DataContractError("Data issue")
3201:     except ContractViolationError as e:
3202:         errors.append(("data", type(e).__name__))
3203:
3204:     try:
3205:         raise SystemContractError("System issue")
3206:     except ContractViolationError as e:
3207:         errors.append(("system", type(e).__name__))
3208:
3209:     assert errors[0] == ("data", "DataContractError")
3210:     assert errors[1] == ("system", "SystemContractError")
3211:     print("    \u2193\233 Data and system errors are distinguishable")
3212:
3213:     print("\n" + "="*60)
3214:     print("Domain exceptions doctest OK - All tests passed!")
3215:     print("="*60)
3216:
3217:
3218:
3219: =====
3220: FILE: src/farfan_pipeline/utils/enhanced_contracts.py
3221: =====
3222:
3223: """
3224: Enhanced Contract System with Pydantic - Production Grade
3225: =====
3226:
3227: Strict contract definitions with cryptographic verification, deterministic
3228: execution guarantees, and comprehensive validation.
3229:
3230: Features:
3231: - Static typing with Pydantic BaseModel
3232: - Schema versioning for backward compatibility
3233: - Cryptographic content digests (SHA-256)
3234: - UTC timestamps (ISO-8601)
3235: - Domain-specific exceptions
3236: - Structured JSON logging
3237: - Flow compatibility validation
3238:
3239: Author: Policy Analytics Research Unit
3240: Version: 2.0.0
3241: License: Proprietary
3242: """
3243:
3244: from __future__ import annotations
3245:
3246: import hashlib
3247: import json
3248: import logging
```

```
3249: import uuid
3250: from datetime import datetime, timezone
3251: from typing import Any
3252:
3253: from pydantic import BaseModel, ConfigDict, Field, field_validator
3254: from farfan_pipeline.core.parameters import ParameterLoaderV2
3255: from farfan_pipeline.core.calibration.decorators import calibrated_method
3256:
3257: # =====
3258: # DOMAIN-SPECIFIC EXCEPTIONS
3259: # =====
3260:
3261: class ContractValidationError(Exception):
3262:     """Raised when contract validation fails."""
3263:
3264:     def __init__(self, message: str, field: str | None = None, event_id: str | None = None) -> None:
3265:         self.field = field
3266:         self.event_id = event_id or str(uuid.uuid4())
3267:         super().__init__(f"[{self.event_id}] {message}")
3268:
3269:
3270: class DataIntegrityError(Exception):
3271:     """Raised when data integrity checks fail (e.g., hash mismatch)."""
3272:
3273:     def __init__(self, message: str, expected: str | None = None, got: str | None = None, event_id: str | None = None) -> None:
3274:         self.expected = expected
3275:         self.got = got
3276:         self.event_id = event_id or str(uuid.uuid4())
3277:         super().__init__(f"[{self.event_id}] {message}")
3278:
3279:
3280: class SystemConfigError(Exception):
3281:     """Raised when system configuration is invalid."""
3282:
3283:     def __init__(self, message: str, config_key: str | None = None, event_id: str | None = None) -> None:
3284:         self.config_key = config_key
3285:         self.event_id = event_id or str(uuid.uuid4())
3286:         super().__init__(f"[{self.event_id}] {message}")
3287:
3288:
3289: class FlowCompatibilityError(Exception):
3290:     """Raised when data flow between components is incompatible."""
3291:
3292:     def __init__(self, message: str, producer: str | None = None, consumer: str | None = None, event_id: str | None = None) -> None:
3293:         self.producer = producer
3294:         self.consumer = consumer
3295:         self.event_id = event_id or str(uuid.uuid4())
3296:         super().__init__(f"[{self.event_id}] {message}")
3297:
3298:
3299: # =====
3300: # UTILITY FUNCTIONS FOR DETERMINISM AND VALIDATION
3301: # =====
3302:
3303: def compute_content_digest(content: str | bytes | dict[str, Any]) -> str:
3304:     """
```

```
3305:     Compute SHA-256 digest of content in a deterministic way.
3306:
3307:     Args:
3308:         content: String, bytes, or dict to hash
3309:
3310:     Returns:
3311:         Hexadecimal SHA-256 digest
3312:
3313:     Examples:
3314:         >>> digest = compute_content_digest("test")
3315:         >>> len(digest)
3316:         64
3317:         >>> digest == compute_content_digest("test")  # Deterministic
3318:         True
3319:     """
3320:     if isinstance(content, dict):
3321:         # Sort keys for deterministic JSON
3322:         content_str = json.dumps(content, sort_keys=True, ensure_ascii=True)
3323:         content_bytes = content_str.encode('utf-8')
3324:     elif isinstance(content, str):
3325:         content_bytes = content.encode('utf-8')
3326:     elif isinstance(content, bytes):
3327:         content_bytes = content
3328:     else:
3329:         raise ContractValidationError(
3330:             f"Cannot compute digest for type {type(content).__name__}",
3331:             field="content"
3332:         )
3333:
3334:     return hashlib.sha256(content_bytes).hexdigest()
3335:
3336:
3337: def utc_now_iso() -> str:
3338:     """
3339:     Get current UTC timestamp in ISO-8601 format.
3340:
3341:     Returns:
3342:         ISO-8601 timestamp string (UTC timezone)
3343:
3344:     Examples:
3345:         >>> ts = utc_now_iso()
3346:         >>> 'T' in ts and 'Z' in ts
3347:         True
3348:     """
3349:     return datetime.now(timezone.utc).isoformat().replace('+00:00', 'Z')
3350:
3351:
3352: # =====
3353: # BASE CONTRACT MODEL
3354: # =====
3355:
3356: class BaseContract(BaseModel):
3357:     """
3358:     Base contract model with common fields for all contracts.
3359:
3360:     All contracts must include:
```

```
3361:     - schema_version: Semantic version for contract evolution
3362:     - timestamp_utc: ISO-8601 UTC timestamp
3363:     - correlation_id: UUID for request tracing
3364:     """
3365:
3366:     model_config = ConfigDict(
3367:         frozen=True, # Immutable for safety
3368:         extra='forbid', # Reject unknown fields
3369:         validate_assignment=True,
3370:         str_strip whitespace=True,
3371:     )
3372:
3373:     schema_version: str = Field(
3374:         default="2.0.0",
3375:         description="Contract schema version (semantic versioning)",
3376:         pattern=r"^\d+\.\d+\.\d+$"
3377:     )
3378:
3379:     timestamp_utc: str = Field(
3380:         default_factory=utc_now_iso,
3381:         description="UTC timestamp in ISO-8601 format"
3382:     )
3383:
3384:     correlation_id: str = Field(
3385:         default_factory=lambda: str(uuid.uuid4()),
3386:         description="UUID for request correlation and tracing"
3387:     )
3388:
3389:     @field_validator('timestamp_utc')
3390:     @classmethod
3391:     def validate_timestamp(cls, v: str) -> str:
3392:         """Validate timestamp is ISO-8601 format and UTC."""
3393:         try:
3394:             dt = datetime.fromisoformat(v.replace('Z', '+00:00'))
3395:             # Ensure UTC
3396:             if dt.tzinfo is None or dt.utcoffset() != timezone.utc.utcoffset(None):
3397:                 raise ValueError("Timestamp must be UTC")
3398:             return v
3399:         except (ValueError, AttributeError) as e:
3400:             raise ContractValidationException(
3401:                 f"Invalid ISO-8601 timestamp: {v}",
3402:                 field="timestamp_utc"
3403:             ) from e
3404:
3405:
3406: # =====
3407: # DOCUMENT CONTRACTS - V2
3408: # =====
3409:
3410: class DocumentMetadataV2(BaseContract):
3411:     """
3412:         Enhanced document metadata with cryptographic verification.
3413:
3414:         Attributes:
3415:             file_path: Absolute path to document
3416:             file_name: Document filename
```

```

3417:     num_pages: Number of pages
3418:     file_size_bytes: File size in bytes
3419:     content_digest: SHA-256 hash of file content
3420:     policy_unit_id: Unique identifier for policy unit
3421:     encoding: Character encoding (default: utf-8)
3422: """
3423:
3424:     file_path: str = Field(..., description="Absolute path to document")
3425:     file_name: str = Field(..., description="Document filename", min_length=1)
3426:     num_pages: int = Field(..., description="Number of pages", ge=1)
3427:     file_size_bytes: int = Field(..., description="File size in bytes", ge=0)
3428:     content_digest: str = Field(..., description="SHA-256 hash of content", pattern=r"^[a-f0-9]{64}$")
3429:     policy_unit_id: str = Field(..., description="Unique policy unit identifier")
3430:     encoding: str = Field(default="utf-8", description="Character encoding")
3431:
3432: # Optional metadata
3433: pdf_metadata: dict[str, Any] | None = Field(default=None, description="PDF metadata dictionary")
3434: author: str | None = Field(default=None, description="Document author")
3435: title: str | None = Field(default=None, description="Document title")
3436: creation_date: str | None = Field(default=None, description="Document creation date")
3437:
3438:
3439: class ProcessedTextV2(BaseContract):
3440: """
3441: Enhanced processed text with input/output validation.
3442:
3443: Attributes:
3444:     raw_text: Original unprocessed text
3445:     normalized_text: Normalized/cleaned text
3446:     language: Detected language code
3447:     input_digest: SHA-256 of raw_text input
3448:     output_digest: SHA-256 of normalized_text output
3449:     policy_unit_id: Policy unit identifier
3450:     processing_latency_ms: Processing time in milliseconds
3451: """
3452:
3453:     raw_text: str = Field(..., description="Original unprocessed text", min_length=1)
3454:     normalized_text: str = Field(..., description="Normalized/cleaned text", min_length=1)
3455:     language: str = Field(..., description="ISO 639-1 language code", pattern=r"^[a-z]{2}$")
3456:     input_digest: str = Field(..., description="SHA-256 of raw_text", pattern=r"^[a-f0-9]{64}$")
3457:     output_digest: str = Field(..., description="SHA-256 of normalized_text", pattern=r"^[a-f0-9]{64}$")
3458:     policy_unit_id: str = Field(..., description="Policy unit identifier")
3459:     processing_latency_ms: float = Field(..., description="Processing latency in ms", ge=ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L236_89", 0.0))
3460:
3461: # Optional fields
3462: sentences: list[str] | None = Field(default=None, description="Sentence segmentation")
3463: sections: list[dict[str, Any]] | None = Field(default=None, description="Document sections")
3464: payload_size_bytes: int | None = Field(default=None, description="Payload size", ge=0)
3465:
3466: @field_validator('input_digest')
3467: @classmethod
3468: def validate_input_digest(cls, v: str, info) -> str:
3469:     """Verify input digest matches raw_text if available."""
3470:     # This is validated post-construction
3471:     return v

```

```
3472:  
3473:  
3474: # =====  
3475: # ANALYSIS CONTRACTS - V2  
3476: # =====  
3477:  
3478: class AnalysisInputV2(BaseContract):  
3479:     """  
3480:         Enhanced analysis input with cryptographic verification.  
3481:  
3482:         Attributes:  
3483:             text: Input text to analyze  
3484:             document_id: Unique document identifier  
3485:             policy_unit_id: Policy unit identifier  
3486:             input_digest: SHA-256 of input text  
3487:             payload_size_bytes: Size of input payload  
3488:     """  
3489:  
3490:     text: str = Field(..., description="Input text to analyze", min_length=1)  
3491:     document_id: str = Field(..., description="Unique document identifier")  
3492:     policy_unit_id: str = Field(..., description="Policy unit identifier")  
3493:     input_digest: str = Field(  
3494:         ...,  
3495:         description="SHA-256 hash of input text",  
3496:         pattern=r"^[a-f0-9]{64}$"  
3497:     )  
3498:     payload_size_bytes: int = Field(..., description="Payload size in bytes", ge=0)  
3499:  
3500:     # Optional context  
3501:     metadata: dict[str, Any] | None = Field(default=None, description="Additional metadata")  
3502:     context: dict[str, Any] | None = Field(default=None, description="Execution context")  
3503:     sentences: list[str] | None = Field(default=None, description="Pre-segmented sentences")  
3504:  
3505:     @classmethod  
3506:     def create_from_text(  
3507:         cls,  
3508:             text: str,  
3509:             document_id: str,  
3510:             policy_unit_id: str,  
3511:             **kwargs: Any  
3512:     ) -> AnalysisInputV2:  
3513:         """  
3514:             Factory method to create AnalysisInputV2 with auto-computed digest.  
3515:  
3516:             Args:  
3517:                 text: Input text  
3518:                 document_id: Document ID  
3519:                 policy_unit_id: Policy unit ID  
3520:                 **kwargs: Additional optional fields  
3521:  
3522:             Returns:  
3523:                 Validated AnalysisInputV2 instance  
3524:         """  
3525:         input_digest = compute_content_digest(text)  
3526:         payload_size_bytes = len(text.encode('utf-8'))  
3527:
```

```
3528:         return cls(
3529:             text=text,
3530:             document_id=document_id,
3531:             policy_unit_id=policy_unit_id,
3532:             input_digest=input_digest,
3533:             payload_size_bytes=payload_size_bytes,
3534:             **kwargs
3535:         )
3536:
3537:
3538: class AnalysisOutputV2(BaseContract):
3539:     """
3540:         Enhanced analysis output with confidence bounds and validation.
3541:
3542:     Attributes:
3543:         dimension: Analysis dimension
3544:         category: Result category
3545:         confidence: Confidence score [ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L322_38", 0,
0), ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L322_43", 1.0)]
3546:         matches: Evidence matches
3547:         output_digest: SHA-256 of output content
3548:         policy_unit_id: Policy unit identifier
3549:         processing_latency_ms: Processing time in milliseconds
3550:     """
3551:
3552:     dimension: str = Field(..., description="Analysis dimension", min_length=1)
3553:     category: str = Field(..., description="Result category", min_length=1)
3554:     confidence: float = Field(..., description="Confidence score", ge=ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L331_70", 0.0), le=ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L331_78", 1.0))
3555:     matches: list[str] = Field(..., description="Evidence matches")
3556:     output_digest: str = Field(..., description="SHA-256 of output", pattern=r"^[a-f0-9]{64}$")
3557:     policy_unit_id: str = Field(..., description="Policy unit identifier")
3558:     processing_latency_ms: float = Field(..., description="Processing latency in ms", ge=ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.FlowCompatibilityError.__init__", "auto_param_L335_89", 0.0))
3559:
3560:     # Optional fields
3561:     positions: list[int] | None = Field(default=None, description="Match positions")
3562:     evidence: list[str] | None = Field(default=None, description="Supporting evidence")
3563:     warnings: list[str] | None = Field(default=None, description="Validation warnings")
3564:     payload_size_bytes: int | None = Field(default=None, description="Output payload size", ge=0)
3565:
3566:     @field_validator('confidence')
3567:     @classmethod
3568:     def validate_confidence_numerical_stability(cls, v: float) -> float:
3569:         """Ensure confidence is numerically stable and within bounds."""
3570:         if not (0.0 <= v <= 1.0):
3571:             raise ContractValidationError(
3572:                 f"Confidence must be in [0.0, 1.0], got {v}",
3573:                 field="confidence"
3574:             )
3575:         # Round to avoid floating point precision issues
3576:         return round(v, 6)
3577:
3578:
3579: # =====
3580: # EXECUTION CONTRACTS - V2
```

```
3581: # =====
3582:
3583: class ExecutionContextV2(BaseContract):
3584:     """
3585:         Enhanced execution context with full observability.
3586:
3587:     Attributes:
3588:         class_name: Executor class name
3589:         method_name: Method being executed
3590:         document_id: Document identifier
3591:         policy_unit_id: Policy unit identifier
3592:         execution_id: Unique execution identifier
3593:         parent_correlation_id: Parent request correlation ID
3594:     """
3595:
3596:     class_name: str = Field(..., description="Executor class name", min_length=1)
3597:     method_name: str = Field(..., description="Method being executed", min_length=1)
3598:     document_id: str = Field(..., description="Document identifier")
3599:     policy_unit_id: str = Field(..., description="Policy unit identifier")
3600:     execution_id: str = Field(
3601:         default_factory=lambda: str(uuid.uuid4())),
3602:         description="Unique execution identifier"
3603:     )
3604:     parent_correlation_id: str | None = Field(
3605:         default=None,
3606:         description="Parent correlation ID for nested calls"
3607:     )
3608:
3609:     # Optional context
3610:     raw_text: str | None = Field(default=None, description="Raw input text")
3611:     text: str | None = Field(default=None, description="Processed text")
3612:     metadata: dict[str, Any] | None = Field(default=None, description="Metadata")
3613:     tables: dict[str, Any] | None = Field(default=None, description="Extracted tables")
3614:     sentences: list[str] | None = Field(default=None, description="Sentences")
3615:
3616:
3617: # =====
3618: # STRUCTURED LOGGING HELPER
3619: # =====
3620:
3621: class StructuredLogger:
3622:     """
3623:         Structured JSON logger for observability.
3624:
3625:     Logs include:
3626:         - correlation_id for tracing
3627:         - latencies per operation
3628:         - payload sizes
3629:         - cryptographic fingerprints
3630:         - NO PII
3631:     """
3632:
3633:     def __init__(self, name: str) -> None:
3634:         """Initialize logger with name."""
3635:         self.logger = logging.getLogger(name)
3636:         self.logger.setLevel(logging.INFO)
```

```
3637:  
3638:     def log_contract_validation(  
3639:         self,  
3640:         contract_type: str,  
3641:         correlation_id: str,  
3642:         success: bool,  
3643:         latency_ms: float,  
3644:         payload_size_bytes: int = 0,  
3645:         content_digest: str | None = None,  
3646:         error: str | None = None  
3647:     ) -> None:  
3648:         """Log contract validation event."""  
3649:         log_entry = {  
3650:             "event": "contract_validation",  
3651:             "contract_type": contract_type,  
3652:             "correlation_id": correlation_id,  
3653:             "success": success,  
3654:             "latency_ms": round(latency_ms, 3),  
3655:             "payload_size_bytes": payload_size_bytes,  
3656:             "timestamp_utc": utc_now_iso(),  
3657:         }  
3658:  
3659:         if content_digest:  
3660:             log_entry["content_digest"] = content_digest  
3661:  
3662:         if error:  
3663:             log_entry["error"] = error  
3664:  
3665:         self.logger.info(json.dumps(log_entry, sort_keys=True))  
3666:  
3667:     def log_execution(  
3668:         self,  
3669:         operation: str,  
3670:         correlation_id: str,  
3671:         success: bool,  
3672:         latency_ms: float,  
3673:         **kwargs: Any  
3674:     ) -> None:  
3675:         """Log execution event with additional context."""  
3676:         log_entry = {  
3677:             "event": "execution",  
3678:             "operation": operation,  
3679:             "correlation_id": correlation_id,  
3680:             "success": success,  
3681:             "latency_ms": round(latency_ms, 3),  
3682:             "timestamp_utc": utc_now_iso(),  
3683:         }  
3684:         log_entry.update(kwargs)  
3685:  
3686:         self.logger.info(json.dumps(log_entry, sort_keys=True))  
3687:  
3688:  
3689: # =====  
3690: # IN-SCRIPT TESTS  
3691: # =====
```

```
3693: if __name__ == "__main__":
3694:     import doctest
3695:
3696:     # Run doctests
3697:     print("Running doctests...")
3698:     doctest.testmod(verbose=True)
3699:
3700:     # Contract validation examples
3701:     print("\n" + "="*60)
3702:     print("Contract Validation Examples")
3703:     print("="*60)
3704:
3705:     # Example 1: Document metadata
3706:     print("\n1. DocumentMetadataV2 validation:")
3707:     doc_meta = DocumentMetadataV2(
3708:         file_path="/path/to/document.pdf",
3709:         file_name="document.pdf",
3710:         num_pages=10,
3711:         file_size_bytes=1024000,
3712:         content_digest="a" * 64, # Valid SHA-256 hex
3713:         policy_unit_id="PDM-001"
3714:     )
3715:     print(f"    \u234\u223 Valid: correlation_id={doc_meta.correlation_id[:8]}...")
3716:
3717:     # Example 2: Analysis input with auto-digest
3718:     print("\n2. AnalysisInputV2 with auto-computed digest:")
3719:     analysis_input = AnalysisInputV2.create_from_text(
3720:         text="Sample policy text for analysis",
3721:         document_id="DOC-123",
3722:         policy_unit_id="PDM-001"
3723:     )
3724:     print(f"    \u234\u223 Valid: input_digest={analysis_input.input_digest[:16]}...")
3725:     print(f"    \u234\u223 Payload size: {analysis_input.payload_size_bytes} bytes")
3726:
3727:     # Example 3: Analysis output with confidence validation
3728:     print("\n3. AnalysisOutputV2 with confidence bounds:")
3729:     analysis_output = AnalysisOutputV2(
3730:         dimension="Dimension1",
3731:         category="CategoryA",
3732:         confidence=ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.StructuredLogger.__init__", "auto_param_L509_19", 0.85), # Must be in [ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.StructuredLogger.__init__", "auto_param_L509_40", 0.0), ParameterLoaderV2.get("farfan_core.utils.enhanced_contracts.StructuredLogger.__init__", "auto_param_L509_45", 1.0)]
3733:         matches=["evidence1", "evidence2"],
3734:         output_digest="b" * 64,
3735:         policy_unit_id="PDM-001",
3736:         processing_latency_ms=123.456
3737:     )
3738:     print(f"    \u234\u223 Valid: confidence={analysis_output.confidence}")
3739:
3740:     # Example 4: Structured logging
3741:     print("\n4. Structured logging example:")
3742:     logger = StructuredLogger("test_logger")
3743:     logger.log_contract_validation(
3744:         contract_type="AnalysisInputV2",
3745:         correlation_id=analysis_input.correlation_id,
3746:         success=True,
```

```
3747:     latency_ms=5.2,
3748:     payload_size_bytes=analysis_input.payload_size_bytes,
3749:     content_digest=analysis_input.input_digest
3750:   )
3751:   print("  \u2192 JSON log emitted to logger")
3752:
3753: # Example 5: Exception handling
3754: print("\n5. Domain-specific exceptions:")
3755: try:
3756:   raise ContractValidationError("Invalid field", field="test_field")
3757: except ContractValidationError as e:
3758:   print(f"  \u2192 ContractValidationError: {e}")
3759:   print(f"  \u2192 Event ID: {e.event_id}")
3760:
3761: print("\n" + "="*60)
3762: print("All validation examples passed!")
3763: print("="*60)
3764:
3765:
3766:
3767: =====
3768: FILE: src/farfan_pipeline/utils/evidence_registry.py
3769: =====
3770:
3771: """Append-only evidence registry with cryptographic hashing.
3772:
3773: This module implements a small ledger that stores evidence entries produced by
3774: analysis components. Each entry links to the previous one through a SHA-256
3775: hash, producing an immutable chain that can be verified for tampering.
3776: """
3777: from __future__ import annotations
3778:
3779: import hashlib
3780: import json
3781: from dataclasses import asdict, dataclass
3782: from datetime import datetime
3783: from pathlib import Path
3784: from typing import TYPE_CHECKING, Any
3785: from farfan_pipeline.core.calibration.decorators import calibrated_method
3786:
3787: if TYPE_CHECKING:
3788:   from collections.abc import Iterable
3789:
3790: def _canonical_json(payload: dict[str, Any]) -> str:
3791:   """Return a canonical JSON representation with sorted keys."""
3792:   return json.dumps(payload, ensure_ascii=False, sort_keys=True, separators=(", ", ":" ))
3793:
3794: @dataclass(frozen=True)
3795: class EvidenceRecord:
3796:   """Single append-only evidence entry."""
3797:
3798:   index: int
3799:   timestamp: str
3800:   method_name: str
3801:   evidence: list[str]
3802:   metadata: dict[str, Any]
```

```
3803:     previous_hash: str
3804:     entry_hash: str
3805:
3806:     @staticmethod
3807:     def create(
3808:         index: int,
3809:         method_name: str,
3810:         evidence: Iterable[str],
3811:         metadata: dict[str, Any] | None,
3812:         previous_hash: str,
3813:         timestamp: datetime | None = None,
3814:     ) -> EvidenceRecord:
3815:         """Build a new evidence record and compute its hash."""
3816:         ts = (timestamp or datetime.utcnow()).isoformat() + "Z"
3817:         metadata_dict = dict(metadata or {})
3818:         evidence_list = list(evidence)
3819:
3820:         payload = {
3821:             "index": index,
3822:             "timestamp": ts,
3823:             "method_name": method_name,
3824:             "evidence": evidence_list,
3825:             "metadata": metadata_dict,
3826:             "previous_hash": previous_hash,
3827:         }
3828:         digest = hashlib.sha256(_canonical_json(payload).encode("utf-8")).hexdigest()
3829:         return EvidenceRecord(
3830:             index=index,
3831:             timestamp=ts,
3832:             method_name=method_name,
3833:             evidence=evidence_list,
3834:             metadata=metadata_dict,
3835:             previous_hash=previous_hash,
3836:             entry_hash=digest,
3837:         )
3838:
3839: class EvidenceRegistry:
3840:     """Append-only registry that persists evidence records to disk."""
3841:
3842:     def __init__(self, storage_path: Path | None = None, auto_load: bool = True) -> None:
3843:         self.storage_path = storage_path or Path(".evidence_registry.json")
3844:         self._records: list[EvidenceRecord] = []
3845:         if auto_load and self.storage_path.exists():
3846:             self._records = self._load_records(self.storage_path)
3847:
3848:     @property
3849:     @calibrated_method("farfan_core.utils.evidence_registry.EvidenceRegistry.records")
3850:     def records(self) -> tuple[EvidenceRecord, ...]:
3851:         """Expose records as an immutable tuple."""
3852:         return tuple(self._records)
3853:
3854:     def append(
3855:         self,
3856:         method_name: str,
3857:         evidence: Iterable[str],
3858:         metadata: dict[str, Any] | None = None,
```

```
3859:         monolith_hash: str | None = None,
3860:     ) -> EvidenceRecord:
3861:         """Append a new evidence record to the registry.
3862:
3863:         Args:
3864:             method_name: Name of the method producing evidence
3865:             evidence: Evidence strings
3866:             metadata: Additional metadata dictionary
3867:             monolith_hash: SHA-256 hash of questionnaire_monolith.json (recommended)
3868:
3869:             ARCHITECTURAL NOTE: Including monolith_hash ensures evidence is
3870:             traceable to the specific questionnaire version that generated it.
3871:             Use factory.compute_monolith_hash() to generate this value.
3872:
3873:             previous_hash = self._records[-1].entry_hash if self._records else "GENESIS"
3874:
3875:             # Merge monolith_hash into metadata if provided
3876:             enriched_metadata = dict(metadata or {})
3877:             if monolith_hash is not None:
3878:                 enriched_metadata['monolith_hash'] = monolith_hash
3879:
3880:             record = EvidenceRecord.create(
3881:                 index=len(self._records),
3882:                 method_name=method_name,
3883:                 evidence=evidence,
3884:                 metadata=enriched_metadata,
3885:                 previous_hash=previous_hash,
3886:             )
3887:             self._records.append(record)
3888:             return record
3889:
3890: # -----
3891: # Persistence
3892: # -----
3893: @calibrated_method("farfan_core.utils.evidence_registry.EvidenceRegistry.save")
3894: def save(self) -> None:
3895:     """Persist the registry to disk."""
3896:     payload = [_serialize_record(record) for record in self._records]
3897:     self.storage_path.write_text(
3898:         json.dumps(payload, indent=2, ensure_ascii=False),
3899:         encoding="utf-8",
3900:     )
3901:
3902: @calibrated_method("farfan_core.utils.evidence_registry.EvidenceRegistry._load_records")
3903: def _load_records(self, path: Path) -> list[EvidenceRecord]:
3904:     try:
3905:         data = json.loads(path.read_text(encoding="utf-8"))
3906:     except json.JSONDecodeError as exc:
3907:         raise ValueError(f"Evidence registry at {path} is not valid JSON: {exc}") from exc
3908:
3909:     if not isinstance(data, list):
3910:         raise ValueError("Evidence registry payload must be a list")
3911:
3912:     records: list[EvidenceRecord] = []
3913:     for index, raw in enumerate(data):
3914:         if not isinstance(raw, dict):
```

```
3915:             raise ValueError("Evidence record must be a JSON object")
3916:     expected_index = raw.get("index")
3917:     if expected_index != index:
3918:         raise ValueError(
3919:             f"Evidence record index mismatch at position {index}: found {expected_index}"
3920:         )
3921:     record = EvidenceRecord(
3922:         index=index,
3923:         timestamp=str(raw.get("timestamp")),
3924:         method_name=str(raw.get("method_name")),
3925:         evidence=list(raw.get("evidence", [])),
3926:         metadata=dict(raw.get("metadata", {})),
3927:         previous_hash=str(raw.get("previous_hash")),
3928:         entry_hash=str(raw.get("entry_hash")),
3929:     )
3930:     records.append(record)
3931:
3932:     self._assert_chain(records)
3933:     return records
3934:
3935: # -----
3936: # Verification utilities
3937: # -----
3938: @calibrated_method("farfan_core.utils.evidence_registry.EvidenceRegistry.verify")
3939: def verify(self) -> bool:
3940:     """Verify registry integrity by recomputing all hashes."""
3941:     self._assert_chain(self._records)
3942:     return True
3943:
3944: @staticmethod
3945: def _assert_chain(records: list[EvidenceRecord]) -> None:
3946:     previous_hash = "GENESIS"
3947:     for expected_index, record in enumerate(records):
3948:         if record.index != expected_index:
3949:             raise ValueError(
3950:                 f"Evidence record out of order: expected index {expected_index}, got {record.index}"
3951:             )
3952:         payload = {
3953:             "index": record.index,
3954:             "timestamp": record.timestamp,
3955:             "method_name": record.method_name,
3956:             "evidence": record.evidence,
3957:             "metadata": record.metadata,
3958:             "previous_hash": previous_hash,
3959:         }
3960:         computed = hashlib.sha256(_canonical_json(payload).encode("utf-8")).hexdigest()
3961:         if computed != record.entry_hash:
3962:             raise ValueError(
3963:                 "Evidence record hash mismatch at index "
3964:                 f"{record.index}: expected {computed}, found {record.entry_hash}"
3965:             )
3966:         previous_hash = record.entry_hash
3967:
3968: def _serialize_record(record: EvidenceRecord) -> dict[str, Any]:
3969:     payload = asdict(record)
3970:     payload["evidence"] = list(record.evidence)
```

```
3971:     payload["metadata"] = dict(record.metadata)
3972:     return payload
3973:
3974: __all__ = [
3975:     "EvidenceRecord",
3976:     "EvidenceRegistry",
3977: ]
3978:
3979:
3980:
3981: =====
3982: FILE: src/farfan_pipeline/utils/hash_utils.py
3983: =====
3984:
3985: """
3986: Hash utilities for deterministic content hashing.
3987:
3988: This module provides cryptographic hashing functions used across the pipeline
3989: for content integrity verification and change detection.
3990:
3991: Author: Integration Team
3992: Version: 1.0.0
3993: Python: 3.10+
3994: """
3995:
3996: import hashlib
3997: import json
3998: from typing import Any
3999:
4000:
4001: def compute_hash(data: dict[str, Any]) -> str:
4002:     """
4003:         Compute deterministic SHA-256 hash of dictionary data.
4004:
4005:         This function creates a canonical JSON representation with sorted keys
4006:         and stable separators to ensure identical dictionaries always produce
4007:         the same hash, regardless of key insertion order.
4008:
4009:     Args:
4010:         data: Dictionary to hash
4011:
4012:     Returns:
4013:         Hexadecimal SHA-256 digest (64 characters)
4014:
4015:     Example:
4016:         >>> data = {"b": 2, "a": 1}
4017:         >>> hash1 = compute_hash(data)
4018:         >>> hash2 = compute_hash({"a": 1, "b": 2})
4019:         >>> hash1 == hash2
4020:         True
4021: """
4022: canonical_json = json.dumps(
4023:     data, sort_keys=True, ensure_ascii=True, separators=(", ", ":"), )
4024: )
4025: return hashlib.sha256(canonical_json.encode("utf-8")).hexdigest()
4026:
```

```
4027:  
4028:  
4029: =====  
4030: FILE: src/farfán_pipeline/utils/json_contract_loader.py  
4031: =====  
4032:  
4033: """Utility helpers to load and validate JSON contract documents."""  
4034: from __future__ import annotations  
4035:  
4036: import hashlib  
4037: import json  
4038: from dataclasses import dataclass  
4039: from pathlib import Path  
4040: from typing import TYPE_CHECKING, Union  
4041: from farfán_pipeline.core.calibration.decorators import calibrated_method  
4042:  
4043: if TYPE_CHECKING:  
4044:     from collections.abc import Iterable, Mapping  
4045:  
4046: PathLike = Union[str, Path]  
4047:  
4048: def _canonical_dump(payload: Mapping[str, object]) -> str:  
4049:     return json.dumps(payload, ensure_ascii=False, sort_keys=True, separators=(", ", ":"))  
4050:  
4051: @dataclass(frozen=True)  
4052: class ContractDocument:  
4053:     """Materialized JSON contract with checksum information."""  
4054:  
4055:     path: Path  
4056:     payload: dict[str, object]  
4057:     checksum: str  
4058:  
4059: @dataclass  
4060: class ContractLoadReport:  
4061:     """Result of attempting to load multiple contract documents."""  
4062:  
4063:     documents: dict[str, ContractDocument]  
4064:     errors: list[str]  
4065:  
4066:     @property  
4067:         @calibrated_method("farfán_core.utils.json_contract_loader.ContractLoadReport.is_successful")  
4068:     def is_successful(self) -> bool:  
4069:         return not self.errors  
4070:  
4071:     @calibrated_method("farfán_core.utils.json_contract_loader.ContractLoadReport.summary")  
4072:     def summary(self) -> str:  
4073:         parts = [f"contracts={len(self.documents)}"]  
4074:         if self.errors:  
4075:             parts.append(f"errors={len(self.errors)}")  
4076:         return ", ".join(parts)  
4077:  
4078: class JSONContractLoader:  
4079:     """Load JSON contract files and compute integrity metadata.  
4080:  
4081:     ARCHITECTURAL BOUNDARY: This loader is for generic JSON contracts ONLY.  
4082:     It must NOT be used to load questionnaire_monolith.json directly.  
4083:
```

```
4083:
4084:     For questionnaire access, use:
4085:         - factory.load_questionnaire() for canonical loading (returns CanonicalQuestionnaire)
4086:         - QuestionnaireResourceProvider for pattern extraction
4087:         """
4088:
4089:     def __init__(self, base_path: Path | None = None) -> None:
4090:         self.base_path = base_path or Path(__file__).resolve().parent
4091:
4092:     @calibrated_method("farfan_core.utils.json_contract_loader.JSONContractLoader.load")
4093:     def load(self, paths: Iterable[PathLike]) -> ContractLoadReport:
4094:         documents: dict[str, ContractDocument] = {}
4095:         errors: list[str] = []
4096:         for raw in paths:
4097:             path = self._resolve_path(raw)
4098:             try:
4099:                 payload = self._read_payload(path)
4100:             except (FileNotFoundException, json.JSONDecodeError, ValueError) as exc:
4101:                 errors.append(f"{path}: {exc}")
4102:             continue
4103:
4104:             checksum = hashlib.sha256(_canonical_dump(payload).encode("utf-8")).hexdigest()
4105:             documents[str(path)] = ContractDocument(path=path, payload=payload, checksum=checksum)
4106:         return ContractLoadReport(documents=documents, errors=errors)
4107:
4108:     @calibrated_method("farfan_core.utils.json_contract_loader.JSONContractLoader.load_directory")
4109:     def load_directory(self, relative_directory: PathLike, pattern: str = "*.json") -> ContractLoadReport:
4110:         directory = self._resolve_path(relative_directory)
4111:         if not directory.exists():
4112:             return ContractLoadReport(documents={}, errors=[f"Directory not found: {directory}"])
4113:         if not directory.is_dir():
4114:             return ContractLoadReport(documents={}, errors=[f"Not a directory: {directory}"])
4115:
4116:         paths = sorted(directory.glob(pattern))
4117:         return self.load(paths)
4118:
4119: # -----
4120: # Helpers
4121: # -----
4122: @calibrated_method("farfan_core.utils.json_contract_loader.JSONContractLoader._resolve_path")
4123: def _resolve_path(self, raw: PathLike) -> Path:
4124:     path = Path(raw)
4125:     if not path.is_absolute():
4126:         path = self.base_path / path
4127:     return path
4128:
4129: @staticmethod
4130: def _read_payload(path: Path) -> dict[str, object]:
4131:     # ARCHITECTURAL GUARD: Block unauthorized questionnaire monolith access
4132:     if path.name == "questionnaire_monolith.json":
4133:         raise ValueError(
4134:             "ARCHITECTURAL VIOLATION: questionnaire_monolith.json must ONLY be "
4135:             "loaded via factory.load_questionnaire() which enforces hash verification. "
4136:             "Use factory.load_questionnaire() for canonical loading."
4137:         )
4138:
```

```
4139:     text = path.read_text(encoding="utf-8")
4140:     data = json.loads(text)
4141:     if not isinstance(data, dict):
4142:         raise ValueError("Contract document must be a JSON object")
4143:     return data
4144:
4145: __all__ = [
4146:     "ContractDocument",
4147:     "ContractLoadReport",
4148:     "JSONContractLoader",
4149: ]
4150:
4151:
```