# README.md

# __init__.py

```python
"""Multi-Agent Market Intelligence System"""

from app.agents.base_agent import BaseAgent
from app.agents.researcher_agent import ResearcherAgent
from app.agents.market_agent import MarketAgent
from app.agents.evaluator_agent import EvaluatorAgent
from app.agents.idea_generator_agent import IdeaGeneratorAgent
from app.agents.orchestrator import Orchestrator
from app.agents.schemas import (
    AgentRole,
    TaskStatus,
    ResearchResult,
    MarketAnalysis,
    GeneratedIdea,
    IdeaEvaluation,
    OrchestrationResult,
    AgentTask,
    AgentMessage,
)

__all__ = [
    "BaseAgent",
    "ResearcherAgent",
    "MarketAgent",
    "EvaluatorAgent",
    "IdeaGeneratorAgent",
    "Orchestrator",
    "AgentRole",
    "TaskStatus",
    "ResearchResult",
    "MarketAnalysis",
    "GeneratedIdea",
    "IdeaEvaluation",
    "OrchestrationResult",
    "AgentTask",
    "AgentMessage",
]
```

# base_agent.py

```python
from abc import ABC, abstractmethod
from typing import Any, Dict, Optional
from datetime import datetime
import uuid

from app.llm.client import LLMWrapper
from app.llm.model_router import ModelRouter
from app.agents.schemas import AgentRole, TaskStatus, AgentTask


class BaseAgent(ABC):

    def __init__(
        self,
        role: AgentRole,
        model: Optional[str] = None,
        system_prompt: Optional[str] = None
    ):
        self.role = role
        self.model = model or ModelRouter.get_model_for_task("synthesis")
        self.llm = LLMWrapper(model=self.model)
        self.system_prompt = system_prompt or self._get_default_system_prompt()
        self.task_history: list[AgentTask] = []

    @abstractmethod
    def _get_default_system_prompt(self) -> str:
        pass

    @abstractmethod
    def execute(self, task: AgentTask) -> AgentTask:
        pass

    def create_task(
        self,
        task_type: str,
        input_data: Dict[str, Any]
    ) -> AgentTask:
        task = AgentTask(
            task_id=str(uuid.uuid4()),
            agent_role=self.role,
            task_type=task_type,
            input_data=input_data,
            status=TaskStatus.PENDING,
            created_at=datetime.now()
        )
        return task

    def _update_task_status(
        self,
        task: AgentTask,
        status: TaskStatus,
        result: Optional[Any] = None,
        error: Optional[str] = None
    ) -> AgentTask:
        task.status = status
        if result:
            task.result = result
        if error:
            task.error = error
        if status == TaskStatus.COMPLETED:
            task.completed_at = datetime.now()
        return task

    def query_llm(
        self,
        prompt: str,
        system_prompt: Optional[str] = None,
        temperature: Optional[float] = None
    ) -> str:
        sys_prompt = system_prompt or self.system_prompt
        return self.llm.query(prompt, system_prompt=sys_prompt)

    def log_task(self, task: AgentTask):
        self.task_history.append(task)

    def get_task_history(self) -> list[AgentTask]:
        return self.task_history

    def __repr__(self) -> str:
        return f"{self.__class__.__name__}(role={self.role.value})"
```

```python
import json
from typing import List
from datetime import datetime

from app.agents.base_agent import BaseAgent
from app.agents.schemas import (
    AgentRole, AgentTask, TaskStatus, IdeaEvaluation,
    GeneratedIdea, MarketAnalysis
)
from app.llm.model_router import ModelRouter


class EvaluatorAgent(BaseAgent):

    def __init__(self):
        super().__init__(
            role=AgentRole.EVALUATOR,
            model=ModelRouter.get_model_for_task("extraction")
        )

    def _get_default_system_prompt(self) -> str:
        return """You are a senior startup evaluator and advisor with experience in:
- Product-market fit assessment
- Market opportunity evaluation
- Feasibility analysis
- Risk assessment
- Investment decision-making

Provide honest, critical, and constructive evaluations. Be specific about strengths,
weaknesses, and actionable recommendations."""

    def execute(self, task: AgentTask) -> AgentTask:
        """Execute an evaluation task"""
        try:
            self._update_task_status(task, TaskStatus.IN_PROGRESS)

            idea = task.input_data.get("idea")
            idea_obj = task.input_data.get("idea_obj")  # GeneratedIdea object
            market_analysis = task.input_data.get("market_analysis")  # MarketAnalysis object

            # Extract idea details
            if idea_obj and isinstance(idea_obj, GeneratedIdea):
                idea_text = f"""
Title: {idea_obj.title}
Description: {idea_obj.description}
Target Audience: {idea_obj.target_audience}
Value Proposition: {idea_obj.value_proposition}
Key Features: {', '.join(idea_obj.key_features)}
Market Opportunity: {idea_obj.market_opportunity}
"""
            else:
                idea_text = idea or "No idea provided"

            market_context = ""
            if market_analysis and isinstance(market_analysis, MarketAnalysis):
                market_context = f"""
Market Context:
- Market Size: {market_analysis.market_size}
- Market Maturity: {market_analysis.market_maturity}
- Trends: {', '.join(market_analysis.trends)}
- Opportunities: {', '.join(market_analysis.opportunities)}
- Threats: {', '.join(market_analysis.threats)}
- Competitive Landscape: {market_analysis.competitive_landscape}
"""

            evaluation_prompt = f"""Evaluate the following startup idea:

{idea_text}

{market_context if market_context else ""}

Provide a comprehensive evaluation in JSON format:
{{
    "feasibility_score": 0.0-1.0,
    "market_potential_score": 0.0-1.0,
    "innovation_score": 0.0-1.0,
    "overall_score": 0.0-1.0,
    "strengths": ["strength 1", "strength 2"],
    "weaknesses": ["weakness 1", "weakness 2"],
    "risks": ["risk 1", "risk 2"],
    "recommendations": ["recommendation 1", "recommendation 2"],
    "verdict": "high_potential|medium_potential|low_potential"
}}

Scoring guidelines:
- Feasibility: Can this be built? (technical, resource, time constraints)
- Market Potential: Is there a real market need? (size, growth, willingness to pay)
- Innovation: How novel/unique is this? (differentiation, competitive advantage)
- Overall: Weighted average considering all factors

Be honest and specific."""

            response = self.query_llm(evaluation_prompt)

            try:
                parsed = json.loads(response)
```

```python
                    except json.JSONDecodeError:
                        parsed = self._parse_evaluation_fallback(response)

                overall_score = parsed.get("overall_score")
                if overall_score is None:
                    overall_score = (
                        parsed.get("feasibility_score", 0.5) * 0.3 +
                        parsed.get("market_potential_score", 0.5) * 0.4 +
                        parsed.get("innovation_score", 0.5) * 0.3
                    )

                result = IdeaEvaluation(
                    idea=idea_text,
                    feasibility_score=float(parsed.get("feasibility_score", 0.5)),
                    market_potential_score=float(parsed.get("market_potential_score", 0.5)),
                    innovation_score=float(parsed.get("innovation_score", 0.5)),
                    overall_score=float(overall_score),
                    strengths=parsed.get("strengths", []),
                    weaknesses=parsed.get("weaknesses", []),
                    risks=parsed.get("risks", []),
                    recommendations=parsed.get("recommendations", []),
                    verdict=parsed.get("verdict", "medium_potential")
                )

                self._update_task_status(task, TaskStatus.COMPLETED, result=result)
                self.log_task(task)

        except Exception as e:
            self._update_task_status(
                task,
                TaskStatus.FAILED,
                error=str(e)
            )
            self.log_task(task)

        return task

    def _parse_evaluation_fallback(self, text: str) -> dict:
        """Fallback parser if JSON parsing fails"""
        return {
            "feasibility_score": 0.5,
            "market_potential_score": 0.5,
            "innovation_score": 0.5,
            "overall_score": 0.5,
            "strengths": [],
            "weaknesses": [],
            "risks": [],
            "recommendations": [],
            "verdict": "medium_potential"
        }

    def evaluate_idea(
        self,
        idea: str = None,
        idea_obj: GeneratedIdea = None,
        market_analysis: MarketAnalysis = None
    ) -> IdeaEvaluation:
        """Convenience method for quick idea evaluation"""
        task = self.create_task(
            "idea_evaluation",
            {
                "idea": idea,
                "idea_obj": idea_obj,
                "market_analysis": market_analysis
            }
        )
        task = self.execute(task)
        return task.result
```

```python
 1 | import json
 2 | from typing import List
 3 | from datetime import datetime
 4 |
 5 | from app.agents.base_agent import BaseAgent
 6 | from app.agents.schemas import (
 7 |     AgentRole, AgentTask, TaskStatus, GeneratedIdea,
 8 |     ResearchResult, MarketAnalysis
 9 | )
10 | from app.llm.model_router import ModelRouter
11 |
12 |
13 | class IdeaGeneratorAgent(BaseAgent):
14 |     """Agent responsible for generating startup ideas based on research"""
15 |
16 |     def __init__(self):
17 |         super().__init__(
18 |             role=AgentRole.IDEA_GENERATOR,
19 |             model=ModelRouter.get_model_for_task("synthesis")
20 |         )
21 |
22 |     def _get_default_system_prompt(self) -> str:
23 |         return """You are a creative startup ideation expert with deep knowledge of:
24 | - Emerging technologies and trends
25 | - Market gaps and opportunities
26 | - User pain points and unmet needs
27 | - Successful startup patterns
28 |
29 | Generate innovative, feasible, and market-driven startup ideas. Focus on:
30 | - Solving real problems
31 | - Clear value propositions
32 | - Specific target audiences
33 | - Differentiated features"""
34 |
35 |     def execute(self, task: AgentTask) -> AgentTask:
36 |         """Execute an idea generation task"""
37 |         try:
38 |             self._update_task_status(task, TaskStatus.IN_PROGRESS)
39 |
40 |             topic = task.input_data.get("topic", "")
41 |             num_ideas = task.input_data.get("num_ideas", 3)
42 |             research_result = task.input_data.get("research_result")
43 |             market_analysis = task.input_data.get("market_analysis")
44 |
45 |             # Build context from research and market analysis
46 |             context_parts = []
47 |
48 |             if research_result and isinstance(research_result, ResearchResult):
49 |                 context_parts.append(f"""
50 | Research Findings:
51 | {research_result.summary}
52 |
53 | Key Insights:
54 | {chr(10).join(f"- {finding}" for finding in research_result.key_findings)}
55 | """)
56 |
57 |             if market_analysis and isinstance(market_analysis, MarketAnalysis):
58 |                 context_parts.append(f"""
59 | Market Analysis:
60 | - Market Size: {market_analysis.market_size}
61 | - Maturity: {market_analysis.market_maturity}
62 | - Trends: {', '.join(market_analysis.trends)}
63 | - Opportunities: {', '.join(market_analysis.opportunities)}
64 | - Threats: {', '.join(market_analysis.threats)}
65 | """)
66 |
67 |             context = "\n".join(context_parts) if context_parts else ""
68 |
69 |             # Generate ideas
70 |             generation_prompt = f"""Generate {num_ideas} innovative startup ideas related to: {topic}
71 |
72 | {context if context else "Base ideas on general market knowledge and trends."}
73 |
74 | For each idea, provide JSON format:
75 | {{
76 |     "title": "Short, catchy title",
77 |     "description": "2-3 sentence description",
78 |     "target_audience": "Specific target audience",
79 |     "value_proposition": "Clear value proposition",
80 |     "key_features": ["feature 1", "feature 2", "feature 3"],
81 |     "market_opportunity": "Why this is a good opportunity",
82 |     "inspiration_sources": ["source 1", "source 2"]
83 | }}
84 |
85 | Return as a JSON array of ideas. Make each idea:
86 | - Specific and actionable
87 | - Based on real pain points or opportunities
88 | - Technically feasible
89 | - Differentiated from existing solutions"""
90 |
91 |             response = self.query_llm(generation_prompt, temperature=0.8)
92 |
93 |             # Parse response
94 |             try:
95 |                 parsed = json.loads(response)
96 |                 # Handle both array and single object responses
```

```python
 97 |                    if isinstance(parsed, dict):
 98 |                        parsed = [parsed]
 99 |                except json.JSONDecodeError:
100 |                    # Try to extract ideas from text
101 |                    parsed = self._parse_ideas_fallback(response)
102 |
103 |                # Create GeneratedIdea objects
104 |                ideas = []
105 |                for idea_data in parsed[:num_ideas]:
106 |                    idea = GeneratedIdea(
107 |                        title=idea_data.get("title", "Untitled Idea"),
108 |                        description=idea_data.get("description", ""),
109 |                        target_audience=idea_data.get("target_audience", ""),
110 |                        value_proposition=idea_data.get("value_proposition", ""),
111 |                        key_features=idea_data.get("key_features", []),
112 |                        market_opportunity=idea_data.get("market_opportunity", ""),
113 |                        inspiration_sources=idea_data.get("inspiration_sources", [])
114 |                    )
115 |                    ideas.append(idea)
116 |
117 |                self._update_task_status(task, TaskStatus.COMPLETED, result=ideas)
118 |                self.log_task(task)
119 |
120 |        except Exception as e:
121 |            self._update_task_status(
122 |                task,
123 |                TaskStatus.FAILED,
124 |                error=str(e)
125 |            )
126 |            self.log_task(task)
127 |
128 |        return task
129 |
130 |    def _parse_ideas_fallback(self, text: str) -> List[dict]:
131 |        """Fallback parser if JSON parsing fails"""
132 |        # Simple extraction - look for numbered or bulleted ideas
133 |        ideas = []
134 |        lines = text.split("\n")
135 |        current_idea = {}
136 |
137 |        for line in lines:
138 |            line = line.strip()
139 |            if not line:
140 |                continue
141 |
142 |            if line.startswith(("1.", "2.", "3.", "4.", "5.")) or \
143 |               (line[0].isupper() and len(line) < 100 and ":" not in line):
144 |                if current_idea:
145 |                    ideas.append(current_idea)
146 |                current_idea = {"title": line.lstrip("1234567890. ")}
147 |            elif ":" in line and current_idea:
148 |                key, value = line.split(":", 1)
149 |                key = key.lower().replace(" ", "_")
150 |                current_idea[key] = value.strip()
151 |
152 |        if current_idea:
153 |            ideas.append(current_idea)
154 |
155 |        return ideas if ideas else [{"title": "Generated Idea", "description": text[:200]}]
156 |
157 |    def generate_ideas(
158 |        self,
159 |        topic: str,
160 |        num_ideas: int = 3,
161 |        research_result=None,
162 |        market_analysis=None
163 |    ) -> List[GeneratedIdea]:
164 |        """Convenience method for quick idea generation"""
165 |        task = self.create_task(
166 |            "idea_generation",
167 |            {
168 |                "topic": topic,
169 |                "num_ideas": num_ideas,
170 |                "research_result": research_result,
171 |                "market_analysis": market_analysis
172 |            }
173 |        )
174 |        task = self.execute(task)
175 |        return task.result if task.result else []
176 |
```

```python
import json
from typing import List, Optional
from datetime import datetime

from app.agents.base_agent import BaseAgent
from app.agents.schemas import AgentRole, AgentTask, TaskStatus, MarketAnalysis, ResearchResult
from app.llm.model_router import ModelRouter


class MarketAgent(BaseAgent):

    def __init__(self):
        super().__init__(
            role=AgentRole.MARKET_ANALYST,
            model=ModelRouter.get_model_for_task("extraction")
        )

    def _get_default_system_prompt(self) -> str:
        return """You are a senior market analyst with expertise in startup ecosystems,
technology markets, and emerging trends. Your role is to:
- Analyze market dynamics and trends
- Assess market size and growth potential
- Identify opportunities and threats
- Evaluate competitive landscapes
- Determine market maturity stages
Provide data-driven, objective analysis."""

    def execute(self, task: AgentTask) -> AgentTask:
        """Execute a market analysis task"""
        try:
            self._update_task_status(task, TaskStatus.IN_PROGRESS)

            topic = task.input_data.get("topic", "")
            research_result = task.input_data.get("research_result")

            context = ""
            if research_result and isinstance(research_result, ResearchResult):
                context = f"""
Research Summary:
{research_result.summary}

Key Findings:
{chr(10).join(f"- {finding}" for finding in research_result.key_findings)}

Sources: {', '.join(research_result.sources)}
"""

            analysis_prompt = f"""Analyze the market for: {topic}

{context if context else "No prior research provided. Base your analysis on general market knowledge."}

Provide a comprehensive market analysis in JSON format:
{{
    "market_size": "Estimated market size (e.g., '$X billion', 'growing rapidly', 'unknown')",
    "trends": ["trend 1", "trend 2", "trend 3"],
    "opportunities": ["opportunity 1", "opportunity 2"],
    "threats": ["threat 1", "threat 2"],
    "competitive_landscape": "Description of competition and market structure",
    "market_maturity": "emerging|growing|mature|declining",
    "confidence_score": 0.0-1.0
}}

Be specific and data-driven where possible."""

            response = self.query_llm(analysis_prompt)

            try:
                parsed = json.loads(response)
            except json.JSONDecodeError:
                parsed = self._parse_market_analysis_fallback(response)

            result = MarketAnalysis(
                topic=topic,
                market_size=parsed.get("market_size"),
                trends=parsed.get("trends", []),
                opportunities=parsed.get("opportunities", []),
                threats=parsed.get("threats", []),
                competitive_landscape=parsed.get("competitive_landscape", ""),
                market_maturity=parsed.get("market_maturity", "unknown"),
                confidence_score=float(parsed.get("confidence_score", 0.5))
            )

            self._update_task_status(task, TaskStatus.COMPLETED, result=result)
            self.log_task(task)

        except Exception as e:
            self._update_task_status(
                task,
                TaskStatus.FAILED,
                error=str(e)
            )
            self.log_task(task)

        return task

    def _parse_market_analysis_fallback(self, text: str) -> dict:
```

```python
 97 |         """Fallback parser if JSON parsing fails"""
 98 |         return {
 99 |             "market_size": "Unknown",
100 |             "trends": [],
101 |             "opportunities": [],
102 |             "threats": [],
103 |             "competitive_landscape": text[:500],
104 |             "market_maturity": "unknown",
105 |             "confidence_score": 0.5
106 |         }
107 |
108 |     def analyze_market(
109 |         self,
110 |         topic: str,
111 |         research_result: Optional[ResearchResult] = None
112 |     ) -> MarketAnalysis:
113 |         """Convenience method for quick market analysis"""
114 |         task = self.create_task(
115 |             "market_analysis",
116 |             {
117 |                 "topic": topic,
118 |                 "research_result": research_result
119 |             }
120 |         )
121 |         task = self.execute(task)
122 |         return task.result
123 |
```

# orchestrator.py

```python
import time
from typing import List, Optional, Dict, Any
from datetime import datetime

from app.agents.base_agent import BaseAgent
from app.agents.researcher_agent import ResearcherAgent
from app.agents.market_agent import MarketAgent
from app.agents.evaluator_agent import EvaluatorAgent
from app.agents.idea_generator_agent import IdeaGeneratorAgent
from app.agents.schemas import (
    AgentRole, AgentTask, TaskStatus, OrchestrationResult,
    ResearchResult, MarketAnalysis, GeneratedIdea, IdeaEvaluation
)
from app.vector_db.config import VectorDBConfig


class Orchestrator:
    """Orchestrates the multi-agent system to complete complex tasks"""

    def __init__(self, vector_db_config: VectorDBConfig):
        self.researcher = ResearcherAgent(vector_db_config)
        self.market_agent = MarketAgent()
        self.evaluator = EvaluatorAgent()
        self.idea_generator = IdeaGeneratorAgent()
        self.execution_history: List[OrchestrationResult] = []

    def run_full_pipeline(
        self,
        topic: str,
        generate_ideas: bool = True,
        evaluate_ideas: bool = True,
        num_ideas: int = 5,
        top_k_research: int = 10
    ) -> OrchestrationResult:
        """
        Run the complete pipeline:
        1. Research the topic
        2. Analyze the market
        3. Generate ideas (if requested)
        4. Evaluate ideas (if requested)
        5. Return top ideas

        Args:
            topic: The topic to research and generate ideas for
            generate_ideas: Whether to generate ideas
            evaluate_ideas: Whether to evaluate generated ideas
            num_ideas: Number of ideas to generate
            top_k_research: Number of research documents to retrieve

        Returns:
            OrchestrationResult with all findings
        """
        start_time = time.time()
        result = OrchestrationResult(
            task_id=f"orchestration_{datetime.now().strftime('%Y%m%d_%H%M%S')}",
            status=TaskStatus.IN_PROGRESS
        )

        try:
            print(f"■ Researching: {topic}")
            research_task = self.researcher.create_task(
                "research",
                {"query": topic, "top_k": top_k_research}
            )
            research_task = self.researcher.execute(research_task)

            if research_task.status == TaskStatus.COMPLETED:
                result.research_results = research_task.result
                print(f"■ Research completed: {len(research_task.result.documents)} documents found")
            else:
                print(f"■ Research failed: {research_task.error}")
                result.status = TaskStatus.FAILED
                return result

            print(f"■ Analyzing market: {topic}")
            market_task = self.market_agent.create_task(
                "market_analysis",
                {
                    "topic": topic,
                    "research_result": result.research_results
                }
            )
            market_task = self.market_agent.execute(market_task)

            if market_task.status == TaskStatus.COMPLETED:
                result.market_analysis = market_task.result
                print(f"■ Market analysis completed: {market_task.result.market_maturity} market")
            else:
                print(f"■ Market analysis failed: {market_task.error}")

            if generate_ideas:
                print(f"■ Generating {num_ideas} ideas...")
                idea_task = self.idea_generator.create_task(
                    "idea_generation",
                    {
                        "topic": topic,
```

```python
                                "num_ideas": num_ideas,
                                "research_result": result.research_results,
                                "market_analysis": result.market_analysis
                            }
                        )
                        idea_task = self.idea_generator.execute(idea_task)

                        if idea_task.status == TaskStatus.COMPLETED:
                            result.generated_ideas = idea_task.result
                            print(f"■ Generated {len(idea_task.result)} ideas")
                        else:
                            print(f"■ Idea generation failed: {idea_task.error}")

                if evaluate_ideas and result.generated_ideas:
                    print(f"■ Evaluating {len(result.generated_ideas)} ideas...")
                    evaluations = []

                    for idea in result.generated_ideas:
                        eval_task = self.evaluator.create_task(
                            "idea_evaluation",
                            {
                                "idea_obj": idea,
                                "market_analysis": result.market_analysis
                            }
                        )
                        eval_task = self.evaluator.execute(eval_task)

                        if eval_task.status == TaskStatus.COMPLETED:
                            evaluations.append(eval_task.result)
                        else:
                            print(f"■■ Evaluation failed for idea: {idea.title}")

                    result.evaluations = evaluations
                    print(f"■ Evaluated {len(evaluations)} ideas")

                    if evaluations:
                        result.top_ideas = self._rank_ideas(
                            result.generated_ideas,
                            evaluations
                        )
                        print(f"■ Top {len(result.top_ideas)} ideas selected")

            result.summary = self._generate_summary(result)
            result.execution_time = time.time() - start_time
            result.status = TaskStatus.COMPLETED

            print(f"\n■ Pipeline completed in {result.execution_time:.2f}s")

        except Exception as e:
            result.status = TaskStatus.FAILED
            result.summary = f"Pipeline failed: {str(e)}"
            result.execution_time = time.time() - start_time
            print(f"■ Pipeline failed: {str(e)}")

        self.execution_history.append(result)
        return result

    def _rank_ideas(
        self,
        ideas: List[GeneratedIdea],
        evaluations: List[IdeaEvaluation]
    ) -> List[GeneratedIdea]:
        """Rank ideas based on evaluations and return top ones"""
        if len(ideas) != len(evaluations):
            return ideas[:3]

        idea_eval_pairs = list(zip(ideas, evaluations))
        idea_eval_pairs.sort(key=lambda x: x[1].overall_score, reverse=True)

        return [idea for idea, _ in idea_eval_pairs[:3]]

    def _generate_summary(self, result: OrchestrationResult) -> str:
        """Generate a summary of the orchestration result"""
        summary_parts = []

        if result.research_results:
            summary_parts.append(
                f"Research: Found {len(result.research_results.documents)} relevant documents "
                f"with {len(result.research_results.key_findings)} key findings."
            )

        if result.market_analysis:
            summary_parts.append(
                f"Market: {result.market_analysis.market_maturity} market with "
                f"{len(result.market_analysis.opportunities)} opportunities identified."
            )

        if result.generated_ideas:
            summary_parts.append(
                f"Ideas: Generated {len(result.generated_ideas)} startup ideas."
            )

        if result.evaluations:
            avg_score = sum(e.overall_score for e in result.evaluations) / len(result.evaluations)
            summary_parts.append(
                f"Evaluation: Average score {avg_score:.2f}/1.0 across {len(result.evaluations)} ideas."
            )
```

```
195 |             if result.top_ideas:
196 |                 summary_parts.append(
197 |                     f"Top Ideas: Selected {len(result.top_ideas)} highest-potential ideas."
198 |                 )
199 |
200 |         return " ".join(summary_parts)
201 |
202 |     def get_execution_history(self) -> List[OrchestrationResult]:
203 |         """Get execution history"""
204 |         return self.execution_history
205 |
206 |     def __repr__(self) -> str:
207 |         return f"Orchestrator(agents=[Researcher, Market, Evaluator, IdeaGenerator])"
208 |
```

```python
 1 | import json
 2 | from typing import List, Dict, Any
 3 | from datetime import datetime
 4 |
 5 | from app.agents.base_agent import BaseAgent
 6 | from app.agents.schemas import AgentRole, AgentTask, TaskStatus, ResearchResult
 7 | from app.rag.retriever import RAGRetriever
 8 | from app.vector_db.config import VectorDBConfig
 9 |
10 |
11 | class ResearcherAgent(BaseAgent):
12 |
13 |     def __init__(self, vector_db_config: VectorDBConfig):
14 |         super().__init__(role=AgentRole.RESEARCHER)
15 |         self.retriever = RAGRetriever(vector_db_config)
16 |
17 |     def _get_default_system_prompt(self) -> str:
18 |         return """You are a senior research analyst specializing in startup and technology markets.
19 | Your role is to:
20 | - Extract key insights from research documents
21 | - Identify patterns and trends
22 | - Summarize findings clearly and concisely
23 | - Highlight the most important information"""
24 |
25 |     def execute(self, task: AgentTask) -> AgentTask:
26 |         try:
27 |             self._update_task_status(task, TaskStatus.IN_PROGRESS)
28 |
29 |             query = task.input_data.get("query", "")
30 |             top_k = task.input_data.get("top_k", 10)
31 |
32 |             documents = self.retriever.retrieve(query, top_k=top_k)
33 |
34 |             context = self._format_documents(documents)
35 |
36 |             summary_prompt = f"""Based on the following research documents, provide:
37 | 1. A comprehensive summary (2-3 paragraphs)
38 | 2. Key findings (bullet points)
39 |
40 | Query: {query}
41 |
42 | Documents:
43 | {context}
44 |
45 | Format your response as JSON:
46 | {{
47 |     "summary": "...",
48 |     "key_findings": ["...", "..."]
49 | }}"""
50 |
51 |             response = self.query_llm(summary_prompt)
52 |
53 |             try:
54 |                 parsed = json.loads(response)
55 |                 summary = parsed.get("summary", "")
56 |                 key_findings = parsed.get("key_findings", [])
57 |             except json.JSONDecodeError:
58 |                 summary = response
59 |                 key_findings = self._extract_findings_from_text(response)
60 |
61 |             sources = [doc.metadata.get("source", "unknown") for doc in documents]
62 |             sources = list(set(sources))
63 |
64 |             result = ResearchResult(
65 |                 query=query,
66 |                 documents=[
67 |                     {
68 |                         "id": doc.id,
69 |                         "text": doc.text,
70 |                         "metadata": doc.metadata,
71 |                         "score": doc.score
72 |                     }
73 |                     for doc in documents
74 |                 ],
75 |                 summary=summary,
76 |                 key_findings=key_findings,
77 |                 sources=sources,
78 |                 timestamp=datetime.now()
79 |             )
80 |
81 |             self._update_task_status(task, TaskStatus.COMPLETED, result=result)
82 |             self.log_task(task)
83 |
84 |         except Exception as e:
85 |             self._update_task_status(
86 |                 task,
87 |                 TaskStatus.FAILED,
88 |                 error=str(e)
89 |             )
90 |             self.log_task(task)
91 |
92 |         return task
93 |
94 |     def _format_documents(self, documents: List) -> str:
95 |         """Format documents for LLM context"""
96 |         formatted = []
```

```
 97 |            for i, doc in enumerate(documents, 1):
 98 |                formatted.append(
 99 |                    f"Document {i} (Relevance: {doc.score:.3f}):\n"
100 |                    f"Source: {doc.metadata.get('source', 'unknown')}\n"
101 |                    f"Content: {doc.text[:500]}...\n"
102 |                )
103 |            return "\n".join(formatted)
104 |
105 |    def _extract_findings_from_text(self, text: str) -> List[str]:
106 |        """Extract findings from text if JSON parsing fails"""
107 |        lines = text.split("\n")
108 |        findings = []
109 |        for line in lines:
110 |            line = line.strip()
111 |            if line and (line.startswith("-") or line.startswith("*") or
112 |                    (line[0].isdigit() and "." in line[:3])):
113 |                findings.append(line.lstrip("-* ").split(". ", 1)[-1])
114 |        return findings[:10]
115 |
116 |    def research(self, query: str, top_k: int = 10) -> ResearchResult:
117 |        """Convenience method for quick research"""
118 |        task = self.create_task("research", {"query": query, "top_k": top_k})
119 |        task = self.execute(task)
120 |        return task.result
121 |
```

```python
from dataclasses import dataclass
from typing import List, Dict, Any, Optional
from datetime import datetime
from enum import Enum


class AgentRole(str, Enum):
    RESEARCHER = "researcher"
    MARKET_ANALYST = "market_analyst"
    EVALUATOR = "evaluator"
    IDEA_GENERATOR = "idea_generator"


class TaskStatus(str, Enum):
    PENDING = "pending"
    IN_PROGRESS = "in_progress"
    COMPLETED = "completed"
    FAILED = "failed"


@dataclass
class ResearchResult:
    query: str
    documents: List[Dict[str, Any]]
    summary: str
    key_findings: List[str]
    sources: List[str]
    timestamp: datetime


@dataclass
class MarketAnalysis:
    topic: str
    market_size: Optional[str]
    trends: List[str]
    opportunities: List[str]
    threats: List[str]
    competitive_landscape: str
    market_maturity: str
    confidence_score: float


@dataclass
class IdeaEvaluation:
    idea: str
    feasibility_score: float
    market_potential_score: float
    innovation_score: float
    overall_score: float
    strengths: List[str]
    weaknesses: List[str]
    risks: List[str]
    recommendations: List[str]
    verdict: str


@dataclass
class GeneratedIdea:
    title: str
    description: str
    target_audience: str
    value_proposition: str
    key_features: List[str]
    market_opportunity: str
    inspiration_sources: List[str]


@dataclass
class AgentTask:
    task_id: str
    agent_role: AgentRole
    task_type: str
    input_data: Dict[str, Any]
    status: TaskStatus = TaskStatus.PENDING
    result: Optional[Any] = None
    error: Optional[str] = None
    created_at: datetime = None
    completed_at: Optional[datetime] = None


@dataclass
class AgentMessage:
    from_agent: AgentRole
    to_agent: AgentRole
    message_type: str
    content: Dict[str, Any]
    timestamp: datetime


@dataclass
class OrchestrationResult:
    task_id: str
    research_results: Optional[ResearchResult] = None
    market_analysis: Optional[MarketAnalysis] = None
    generated_ideas: List[GeneratedIdea] = None
    evaluations: List[IdeaEvaluation] = None
```

```
 97 |     top_ideas: List[GeneratedIdea] = None
 98 |     summary: str = ""
 99 |     execution_time: float = 0.0
100 |     status: TaskStatus = TaskStatus.PENDING
101 |
```

# settings.py

```python
1  | from functools import lru_cache
2  | from typing import Literal
3  |
4  | from pydantic_settings import BaseSettings, SettingsConfigDict
5  |
6  |
7  | # ========================
8  | # Model Configuration
9  | # ========================
10 | class ModelConfig(BaseSettings):
11 |     CHAT_MODEL: str = "gpt-4o-mini"
12 |     EMBEDDING_MODEL: str = "text-embedding-3-small"
13 |     TEMPERATURE: float = 0.3
14 |     MAX_TOKENS: int = 1000
15 |
16 |
17 | class ModelRoutingConfig(BaseSettings):
18 |     IDEATION_MODEL: str = "gpt-4o-mini"
19 |     ANALYSIS_MODEL: str = "gpt-4o"
20 |     SUMMARIZATION_MODEL: str = "gpt-4o-mini"
21 |
22 |
23 | # ========================
24 | # Vector Database Config
25 | # ========================
26 | class VectorDBConfig(BaseSettings):
27 |     PROVIDER: Literal["chroma", "pinecone", "qdrant"] = "chroma"
28 |     COLLECTION_NAME: str = "documents"
29 |     PERSIST_DIR: str = "./vector_store"
30 |
31 |
32 | # ========================
33 | # Logging + Observability
34 | # ========================
35 | class LoggingConfig(BaseSettings):
36 |     LOG_LEVEL: Literal["DEBUG", "INFO", "WARNING", "ERROR", "CRITICAL"] = "INFO"
37 |     LOG_FILE: str = "./logs/app.log"
38 |
39 |
40 | class ObservabilityConfig(BaseSettings):
41 |     ENABLE_TRACING: bool = True
42 |     TRACE_LLM_CALLS: bool = True
43 |     TRACE_PIPELINES: bool = True
44 |
45 |
46 | # ========================
47 | # Scraper Config
48 | # ========================
49 | class ScraperConfig(BaseSettings):
50 |     REDDIT_LIMIT: int = 100
51 |     HN_LIMIT: int = 100
52 |     X_LIMIT: int = 100
53 |     REQUEST_TIMEOUT: int = 10
54 |     RETRIES: int = 3
55 |     RATE_LIMIT_PER_MINUTE: int = 30
56 |     BACKOFF_FACTOR: float = 1.5
57 |
58 |
59 | # ========================
60 | # LLM Control
61 | # ========================
62 | class LLMConfig(BaseSettings):
63 |     MAX_RETRIES: int = 3
64 |     TIMEOUT: int = 60
65 |     COST_LIMIT_PER_RUN: float = 5.0
66 |
67 |
68 | # ========================
69 | # Pipeline Config
70 | # ========================
71 | class PipelineConfig(BaseSettings):
72 |     BATCH_SIZE: int = 32
73 |     EMBEDDING_BATCH_SIZE: int = 64
74 |     ENABLE_PARALLEL_STAGES: bool = True
75 |
76 |
77 | # ========================
78 | # Runtime Flags
79 | # ========================
80 | class RuntimeFlags(BaseSettings):
81 |     ENABLE_LLM: bool = True
82 |     ENABLE_SCRAPING: bool = True
83 |     DRY_RUN: bool = False
84 |
85 |
86 | # ========================
87 | # Cache Config
88 | # ========================
89 | class CacheConfig(BaseSettings):
90 |     ENABLE_CACHE: bool = True
91 |     CACHE_DIR: str = "./cache"
92 |
93 |
94 | # ========================
95 | # Worker / Async Config
96 | # ========================
```

```python
 97 | class WorkerConfig(BaseSettings):
 98 |     MAX_CONCURRENT_TASKS: int = 5
 99 |     QUEUE_SIZE: int = 100
100 |
101 |
102 | # =========================
103 | # Paths Config
104 | # =========================
105 | class Paths(BaseSettings):
106 |     RAW_DATA: str = "data/raw"
107 |     PROCESSED_DATA: str = "data/processed"
108 |     REPORTS: str = "reports"
109 |
110 |
111 | # =========================
112 | # Root Application Settings
113 | # =========================
114 | class AppSettings(BaseSettings):
115 |     model_config = SettingsConfigDict(env_file=".env")
116 |
117 |     OPENAI_API_KEY: str
118 |     PINECONE_API_KEY: str | None = None
119 |     ENV: Literal["development", "staging", "production"] = "development"
120 |
121 |     model: ModelConfig = ModelConfig()
122 |     model_routing: ModelRoutingConfig = ModelRoutingConfig()
123 |     vectordb: VectorDBConfig = VectorDBConfig()
124 |
125 |     logging: LoggingConfig = LoggingConfig()
126 |     observability: ObservabilityConfig = ObservabilityConfig()
127 |
128 |     scraper: ScraperConfig = ScraperConfig()
129 |     llm: LLMConfig = LLMConfig()
130 |     pipeline: PipelineConfig = PipelineConfig()
131 |
132 |     runtime_flags: RuntimeFlags = RuntimeFlags()
133 |     cache: CacheConfig = CacheConfig()
134 |     worker: WorkerConfig = WorkerConfig()
135 |     paths: Paths = Paths()
136 |
137 |
138 | # =========================
139 | # Settings Factory (Singleton)
140 | # =========================
141 | @lru_cache
142 | def get_settings() -> AppSettings:
143 |     return AppSettings()
144 |
145 |
146 | # default instance
147 | settings = get_settings()
148 |
```

# hn_scraper.py

```python
import requests
import json
import os
from datetime import datetime
from app.config.settings import get_settings

settings = get_settings()

HN_API_URL = "http://hn.algolia.com/api/v1/search?tags=front_page"

def scrape_hackernews(limit=100):
    all_posts = []

    response = requests.get(HN_API_URL)
    response.raise_for_status()
    data = response.json()

    for post in data.get("hits", [])[:limit]:
        all_posts.append({
            "title": post.get("title"),
            "body": post.get("story_text") or "",
            "score": post.get("points"),
            "comments": post.get("num_comments"),
            "created_at": datetime.fromtimestamp(post.get("created_at_i")).isoformat()
            if post.get("created_at_i") else None,
            "source": "hackernews"
        })

    file_path = os.path.join(settings.paths.RAW_DATA, "hn.json")
    with open(file_path, "w", encoding="utf-8") as f:
        json.dump(all_posts, f, indent=2, ensure_ascii=False)

    print(f"Scraped {len(all_posts)} Hacker News posts -> {file_path}")
    return file_path

if __name__ == "__main__":
    scrape_hackernews(limit=settings.scraper.HN_LIMIT)
```

# reddit_scraper.py

```python
import json
import os
from datetime import datetime
import praw
from app.config.settings import get_settings

settings = get_settings()

SUBREDDITS = ['startups', 'Entrepreneur', 'SaaS', 'SideProject']

def scrape_reddit(limit=100):
    reddit = praw.Reddit(
        client_id=os.getenv("REDDIT_CLIENT_ID"),
        client_secret=os.getenv("REDDIT_CLIENT_SECRET"),
        user_agent=os.getenv("REDDIT_USER_AGENT")
    )

    all_posts = []

    for sub in SUBREDDITS:
        subreddit = reddit.subreddit(sub)
        for post in subreddit.hot(limit=limit):
            all_posts.append({
                "title": post.title,
                "body": post.selftext,
                "score": post.score,
                "comments": post.num_comments,
                "created_at": datetime.fromtimestamp(post.created_utc).isoformat(),
                "source": f"reddit/r/{sub}"
            })

    file_path = os.path.join(settings.paths.RAW_DATA, "reddit.json")
    os.makedirs(os.path.dirname(file_path), exist_ok=True)
    with open(file_path, "w", encoding="utf-8") as f:
        json.dump(all_posts, f, ensure_ascii=False)

    print(f"Scraped {len(all_posts)} Reddit posts -> {file_path}")
    return file_path

if __name__ == "__main__":
    scrape_reddit(limit=settings.scraper.REDDIT_LIMIT)
```

# test_json.py

```python
import json
import os
from app.config.settings import get_settings

settings = get_settings()

def test_json(file_name):
    path = os.path.join(settings.paths.RAW_DATA, file_name)
    try:
        with open(path, "r", encoding="utf-8") as f:
            data = json.load(f)
        print(f"{file_name} is valid JSON with {len(data)} records")
    except Exception as e:
        print(f"{file_name} is invalid JSON! Error: {e}")

if __name__ == "__main__":
    test_json("reddit.json")
    test_json("hn.json")
```

# cache_manager.py

```python
import hashlib
import json
from app.config.settings import get_settings

settings = get_settings()

class CacheManager:
    def __init__(self):
        self.cache = {}

    def _generate_key(self, model: str, messages: list) -> str:
        data = f"{model}:{json.dumps(messages, sort_keys=True)}"
        return hashlib.md5(data.encode()).hexdigest()

    def get(self, model: str, messages: list):
        key = self._generate_key(model, messages)
        return self.cache.get(key)

    def set(self, model: str, messages: list, response: str):
        key = self._generate_key(model, messages)
        if settings.cache.ENABLE_CACHE:
            self.cache[key] = response
```

# client.py

```python
from app.config.settings import get_settings
from app.llm.cache_manager import CacheManager
from app.llm.retry_handler import retry_with_backoff
from app.llm.model_router import ModelRouter
from openai import OpenAI

settings = get_settings()

class LLMWrapper:
    def __init__(self, model: str = None):
        self.model = model or settings.model_routing.IDEATION_MODEL
        self.cache = CacheManager()
        self.client = OpenAI(api_key=settings.OPENAI_API_KEY)

    @retry_with_backoff(retries=settings.llm.MAX_RETRIES, backoff_in_seconds=1)
    def query(self, prompt: str, system_prompt="You are a helpful assistant"):
        cached = self.cache.get(self.model, [{"role": "user", "content": prompt}])
        if cached:
            return cached

        response = self.client.chat.completions.create(
            model=self.model,
            messages=[
                {"role": "system", "content": system_prompt},
                {"role": "user", "content": prompt}
            ],
            temperature=settings.model.TEMPERATURE,
            max_tokens=settings.model.MAX_TOKENS
        )

        answer = response.choices[0].message.content
        self.cache.set(self.model, [{"role": "user", "content": prompt}], answer)
        return answer
```

# model_router.py

```python
from app.config.settings import get_settings

settings = get_settings()

class ModelRouter:
    MODEL_MAP = {
        "synthesis": settings.model_routing.IDEATION_MODEL,
        "extraction": settings.model_routing.ANALYSIS_MODEL,
        "classification": settings.model_routing.CLASSIFICATION_MODEL
        if hasattr(settings.model_routing, "CLASSIFICATION_MODEL")
        else settings.model_routing.ANALYSIS_MODEL
    }

    @classmethod
    def get_model_for_task(cls, task_type: str):
        return cls.MODEL_MAP.get(task_type, settings.model_routing.IDEATION_MODEL)
```

# template.py

```python
MARKET_ANALYSIS_TEMPLATE = """
Analyze the following market data for {sector}:
Data: {data}
Focus on: {key_metrics}
"""

def format_analysis_prompt(sector: str, data: str, key_metrics: str) -> str:
    return MARKET_ANALYSIS_TEMPLATE.format(
        sector=sector,
        data=data,
        key_metrics=key_metrics
    )
```

# retry_handler.py

```python
import time
from functools import wraps

def retry_with_backoff(retries=3, backoff_in_seconds=1):
    def decorator(func):
        @wraps(func)
        def wrapper(*args, **kwargs):
            attempts = 0
            while attempts < retries:
                try:
                    return func(*args, **kwargs)
                except Exception as e:
                    attempts += 1
                    print(f"Error: {e}. Retrying {attempts}/{retries}...")
                    time.sleep(backoff_in_seconds * (2 ** (attempts - 1)))
            return func(*args, **kwargs)
        return wrapper
    return decorator
```

# main.py

```
1
2
```

# context_compressor.py

```python
from typing import List
from app.rag.schemas import RetrievedDocument


class ContextCompressor:
    def __init__(self, max_chars: int = 8000):
        self.max_chars = max_chars

    def compress(self, docs: List[RetrievedDocument]) -> str:
        context = ""

        for doc in docs:
            chunk = f"[Source:{doc.metadata.get('source')}]\n{doc.text}\n\n"

            if len(context) + len(chunk) > self.max_chars:
                break

            context += chunk

        return context.strip()
```

# insight_generator.py

```python
from app.llm.client import LLMWrapper
from app.llm.model_router import ModelRouter
from app.rag.schemas import InsightResult
import json


class InsightGenerator:
    def __init__(self):
        model = ModelRouter.get_model_for_task("synthesis")
        self.llm = LLMWrapper(model=model)

    def generate(self, query: str, context: str) -> InsightResult:

        system_prompt = """
You are a senior startup analyst.
Extract:
- key pain points
- market opportunities
- emerging trends
- strong signals
Return JSON.
"""

        prompt = f"""
QUERY:
{query}

DATA:
{context}

OUTPUT FORMAT:
{{
 "summary": "",
 "pain_points": [],
 "opportunities": [],
 "signals": []
}}
"""

        response = self.llm.query(prompt, system_prompt=system_prompt)

        data = json.loads(response)

        return InsightResult(
            summary=data.get("summary", ""),
            pain_points=data.get("pain_points", []),
            opportunities=data.get("opportunities", []),
            signals=data.get("signals", []),
        )
```

# painpoint_clusterer.py

```python
from collections import defaultdict
from typing import List


class PainPointClusterer:
    def cluster(self, pain_points: List[str]):

        clusters = defaultdict(list)

        for p in pain_points:
            key = p.split()[0]
            clusters[key].append(p)

        return dict(clusters)
```

# retriever.py

```python
from typing import List
from app.vector_db.semantic_search import SemanticSearchEngine
from app.vector_db.config import VectorDBConfig
from app.rag.schemas import RetrievedDocument


class RAGRetriever:
    def __init__(self, config: VectorDBConfig):
        self.engine = SemanticSearchEngine(config)

    def retrieve(self, query: str, top_k: int = 10) -> List[RetrievedDocument]:
        results = self.engine.search(query, top_k=top_k)

        docs = []

        ids = results["ids"][0]
        texts = results["documents"][0]
        metas = results["metadatas"][0]
        scores = results["distances"][0]

        for i in range(len(ids)):
            docs.append(
                RetrievedDocument(
                    id=ids[i],
                    text=texts[i],
                    metadata=metas[i],
                    score=scores[i],
                )
            )

        return docs
```

# schemas.py

```python
from dataclasses import dataclass
from typing import List, Dict, Any


@dataclass
class RetrievedDocument:
    id: str
    text: str
    metadata: Dict[str, Any]
    score: float


@dataclass
class InsightResult:
    summary: str
    pain_points: List[str]
    opportunities: List[str]
    signals: List[str]
```

# trend_detector.py

```python
from collections import Counter
from typing import List
from app.rag.schemas import RetrievedDocument


class TrendDetector:
    def detect(self, docs: List[RetrievedDocument], top_k=5):

        words = []

        for d in docs:
            words.extend(d.text.split())

        counter = Counter(words)

        return counter.most_common(top_k)
```

# cache.py

```python
import json
from pathlib import Path


class EmbeddingCache:
    def __init__(self, path="data/cache/embeddings.json"):
        self.path = Path(path)
        self.path.parent.mkdir(parents=True, exist_ok=True)

        if self.path.exists():
            self.cache = json.loads(self.path.read_text())
        else:
            self.cache = {}

    def get(self, key):
        return self.cache.get(key)

    def set(self, key, value):
        self.cache[key] = value
        self.path.write_text(json.dumps(self.cache))
```

# chroma_client.py

```python
import chromadb
from chromadb.config import Settings
from app.vector_db.config import VectorDBConfig


class ChromaVectorDB:
    def __init__(self, config: VectorDBConfig):
        self.config = config
        self.client = chromadb.Client(Settings(
            chroma_db_impl="duckdb+parquet",
            persist_directory=self.config.persist_directory,
            anonymized_telemetry=False
        ))
        self.collection = self.client.get_or_create_collection(name=self.config.collection_name)

    def add_documents(self, ids: list[str], documents: list[str], embeddings: list[list[float]], metadatas: list[dict] | None = N
        self.collection.add(
            ids=ids,
            documents=documents,
            embeddings=embeddings,
            metadatas=metadatas
        )

    def query(self, query_embedding, top_k=5):
        return self.collection.query(
            query_embeddings=[query_embedding],
            n_results=top_k
        )
```

# config.py

```python
from dataclasses import dataclass

@dataclass
class VectorDBConfig:
    persist_directory: str = "data/vector_db"
    collection_name: str = "startup_ideas"
    embedding_model: str = "text-embedding-3-small"
    batch_size: int = 32
    max_retries: int = 3

```

# deduplicator.py

```python
import hashlib


def hash_text(text: str) -> str:
    return hashlib.sha256(text.encode()).hexdigest()
```

# embedding_service.py

```python
import time
from openai import OpenAI
from app.vector_db.logger import get_logger
from app.config.settings import get_settings

settings = get_settings()
logger = get_logger("embedding_service")


class EmbeddingService:
    def __init__(self, model: str, max_retries=3):
        self.client = OpenAI(api_key=settings.OPENAI_API_KEY)
        self.model = model
        self.max_retries = max_retries

    def embed_batch(self, texts):
        for attempt in range(self.max_retries):
            try:
                response = self.client.embeddings.create(
                    model=self.model,
                    input=texts,
                )
                return [e.embedding for e in response.data]

            except Exception as e:
                wait = 2 ** attempt
                logger.warning(f"Embedding failed, retrying in {wait}s | {e}")
                time.sleep(wait)

        raise RuntimeError("Embedding failed after retries")
```

```python
 1 | import json
 2 | from tqdm import tqdm
 3 |
 4 | from app.vector_db.schemas import Document
 5 | from app.vector_db.deduplicator import hash_text
 6 | from app.vector_db.embedding_service import EmbeddingService
 7 | from app.vector_db.cache import EmbeddingCache
 8 | from app.vector_db.chroma_client import ChromaVectorDB
 9 | from app.vector_db.logger import get_logger
10 | from app.vector_db.config import VectorDBConfig
11 |
12 |
13 | logger = get_logger("ingestion_pipeline")
14 |
15 |
16 | class VectorIngestionPipeline:
17 |     def __init__(self, config: VectorDBConfig):
18 |         self.config = config
19 |         self.db = ChromaVectorDB(config)
20 |         self.embedder = EmbeddingService(
21 |             model=config.embedding_model,
22 |             max_retries=config.max_retries,
23 |         )
24 |         self.cache = EmbeddingCache()
25 |
26 |     def load_json(self, path: str):
27 |         with open(path, "r", encoding="utf-8") as f:
28 |             data = json.load(f)
29 |
30 |         documents = []
31 |
32 |         for item in data:
33 |             text = f"{item.get('title','')} {item.get('body','')}".strip()
34 |
35 |             doc = Document(
36 |                 id=hash_text(text),
37 |                 text=text,
38 |                 metadata={
39 |                     "source": item.get("source"),
40 |                     "score": item.get("score"),
41 |                     "comments": item.get("comments"),
42 |                     "created_at": item.get("created_at"),
43 |                 },
44 |             )
45 |             documents.append(doc)
46 |
47 |         return documents
48 |
49 |     def index(self, documents):
50 |         batch_size = self.config.batch_size
51 |
52 |         for i in tqdm(range(0, len(documents), batch_size)):
53 |             batch = documents[i : i + batch_size]
54 |
55 |             texts = []
56 |             ids = []
57 |             metas = []
58 |             embeddings = []
59 |
60 |             for doc in batch:
61 |                 cached = self.cache.get(doc.id)
62 |
63 |                 if cached:
64 |                     embeddings.append(cached)
65 |                 else:
66 |                     texts.append(doc.text)
67 |
68 |                 ids.append(doc.id)
69 |                 metas.append(doc.metadata)
70 |
71 |             if texts:
72 |                 new_embeddings = self.embedder.embed_batch(texts)
73 |
74 |                 for text, emb in zip(texts, new_embeddings):
75 |                     self.cache.set(hash_text(text), emb)
76 |                     embeddings.append(emb)
77 |
78 |             self.db.add_documents(
79 |                 ids=ids,
80 |                 embeddings=embeddings,
81 |                 documents=[d.text for d in batch],
82 |                 metadatas=metas,
83 |             )
84 |
85 |             logger.info(f"Indexed batch size={len(batch)}")
86 |
```

# logger.py

```python
import logging

def get_logger(name: str):
    logger = logging.getLogger(name)
    logger.setLevel(logging.INFO)

    if not logger.handlers:
        handler = logging.StreamHandler()
        formatter = logging.Formatter(
            '%(asctime)s - %(name)s - %(levelname)s - %(message)s'
        )
        handler.setFormatter(formatter)
        logger.addHandler(handler)

    return logger
```

# schemas.py

```python
from dataclasses import dataclass
from typing import Dict, Any

@dataclass
class Document:
    id: str
    text: str
    metadata: Dict[str, Any]
```

# semantic_search.py

```python
from openai import OpenAI
from app.vector_db.chroma_client import ChromaVectorDB
from app.vector_db.config import VectorDBConfig
from app.config.settings import get_settings

settings = get_settings()


class SemanticSearchEngine:
    def __init__(self, config: VectorDBConfig):
        self.client = OpenAI(api_key=settings.OPENAI_API_KEY)
        self.model = config.embedding_model
        self.db = ChromaVectorDB(config)

    def embed_query(self, query):
        response = self.client.embeddings.create(
            model=self.model,
            input=[query],
        )
        return response.data[0].embedding

    def search(self, query, top_k=5):
        embedding = self.embed_query(query)
        return self.db.query(embedding, top_k)
```

# similarity.py

```
1 | import math
2 |
3 |
4 | def cosine_similarity(a, b):
5 |     dot = sum(x * y for x, y in zip(a, b))
6 |     na = math.sqrt(sum(x * x for x in a))
7 |     nb = math.sqrt(sum(y * y for y in b))
8 |     return dot / (na * nb) if na and nb else 0.0
9 |
```

```python
 1  import os
 2  from reportlab.platypus import SimpleDocTemplate, Paragraph, Spacer, Preformatted, PageBreak
 3  from reportlab.lib.styles import getSampleStyleSheet, ParagraphStyle
 4  from reportlab.lib.pagesizes import A4
 5  from reportlab.lib.units import mm
 6  from reportlab.lib import colors
 7  from pygments import highlight
 8  from pygments.lexers import get_lexer_for_filename, TextLexer
 9  from pygments.formatters import HtmlFormatter
10  from pygments.styles import get_style_by_name
11  from reportlab.lib.styles import StyleSheet1
12  from reportlab.platypus import Flowable
13  from reportlab.lib.enums import TA_CENTER
14  from reportlab.lib.styles import ParagraphStyle
15  from reportlab.lib.styles import getSampleStyleSheet
16  from reportlab.lib.styles import ParagraphStyle
17  from reportlab.platypus import XPreformatted
18
19
20  EXCLUDE_DIRS = {".git", "__pycache__", "venv", ".venv", "node_modules", "build", "dist"}
21  SUPPORTED_EXTENSIONS = (".py", ".js", ".ts", ".go", ".java", ".cpp", ".c", ".json", ".yaml", ".yml", ".md")
22
23
24  def collect_code_files(directory):
25      code_files = []
26      for root, dirs, files in os.walk(directory):
27          dirs[:] = [d for d in dirs if d not in EXCLUDE_DIRS]
28
29          for file in files:
30              if file.endswith(SUPPORTED_EXTENSIONS):
31                  code_files.append(os.path.join(root, file))
32
33      return sorted(code_files)
34
35
36  def highlight_code(code, filename):
37      try:
38          lexer = get_lexer_for_filename(filename)
39      except:
40          lexer = TextLexer()
41
42      formatter = HtmlFormatter(style="monokai", noclasses=True)
43      return highlight(code, lexer, formatter)
44
45
46  def strip_html_tags(text):
47      import re
48      clean = re.compile("<.*?>")
49      return re.sub(clean, "", text)
50
51
52  def create_pdf_from_code(files, output_pdf):
53      doc = SimpleDocTemplate(
54          output_pdf,
55          pagesize=A4,
56          rightMargin=20,
57          leftMargin=20,
58          topMargin=20,
59          bottomMargin=20,
60      )
61
62      styles = getSampleStyleSheet()
63
64      title_style = ParagraphStyle(
65          name="TitleStyle",
66          fontName="Helvetica-Bold",
67          fontSize=12,
68          spaceAfter=8,
69          alignment=TA_CENTER,
70      )
71
72      code_style = ParagraphStyle(
73          name="CodeStyle",
74          fontName="Courier",
75          fontSize=7,
76          leading=8,
77          backColor=None,
78      )
79
80      elements = []
81
82      for file in files:
83          with open(file, "r", encoding="utf-8", errors="ignore") as f:
84              code = f.read()
85
86          elements.append(Paragraph(os.path.basename(file), title_style))
87
88          highlighted = highlight_code(code, file)
89          clean_code = strip_html_tags(highlighted)
90
91          numbered_lines = "\n".join(
92              f"{str(i+1).rjust(4)} | {line}"
93              for i, line in enumerate(clean_code.splitlines())
94          )
95
96          elements.append(XPreformatted(numbered_lines, code_style))
```

```
 97 |             elements.append(PageBreak())
 98 |
 99 |     doc.build(elements)
100 |
101 |
102 | if __name__ == "__main__":
103 |     project_dir = "./"
104 |     output_pdf = "project_code_pro.pdf"
105 |
106 |     code_files = collect_code_files(project_dir)
107 |
108 |     print(f"Found {len(code_files)} files")
109 |     create_pdf_from_code(code_files, output_pdf)
110 |
111 |     print(f"PDF created: {output_pdf}")
112 |
```

# data_pipeline.py

```python
1   | import json
2   | import re
3   | from pathlib import Path
4   | from typing import Iterable, List
5   |
6   | from app.config.settings import get_settings
7   | from app.llm.client import LLMWrapper
8   | from app.llm.model_router import ModelRouter
9   | from app.llm.prompts.template import format_analysis_prompt
10  |
11  | settings = get_settings()
12  |
13  | DEFAULT_SECTOR = "startups & emerging tech"
14  |
15  | DEFAULT_KEY_METRICS = (
16  |     "pain points, unmet needs, traction signals, "
17  |     "pricing hints, and competitive alternatives"
18  | )
19  |
20  | def load_json_records(path: Path) -> List[dict]:
21  |     if not path.exists():
22  |         raise FileNotFoundError(f"File not found: {path}")
23  |
24  |     with path.open("r", encoding="utf-8") as handle:
25  |         return json.load(handle)
26  |
27  | def preprocess_text(title: str | None, body: str | None) -> str:
28  |     combined = " ".join([segment for segment in [title, body] if segment]).strip()
29  |
30  |     combined = combined.lower()
31  |     combined = re.sub(r"http\S+|www\.\S+", "", combined)
32  |     combined = re.sub(r"[^\w\s\-\.,!?]", "", combined)
33  |     combined = re.sub(r"\s+", " ", combined)
34  |
35  |     return combined.strip()
36  |
37  | def iter_batches(items: List[str], batch_size: int) -> Iterable[List[str]]:
38  |     for start in range(0, len(items), batch_size):
39  |         yield items[start:start + batch_size]
40  |
41  | def generate_embeddings(texts: List[str]) -> List[List[float]]:
42  |     if not texts:
43  |         return []
44  |
45  |     from openai import OpenAI
46  |
47  |     client = OpenAI(api_key=settings.OPENAI_API_KEY)
48  |
49  |     embeddings: List[List[float]] = []
50  |
51  |     for batch in iter_batches(texts, batch_size=50):
52  |         response = client.embeddings.create(
53  |             model=settings.model.EMBEDDING_MODEL,
54  |             input=batch
55  |         )
56  |
57  |         ordered = sorted(response.data, key=lambda item: item.index)
58  |         embeddings.extend([item.embedding for item in ordered])
59  |
60  |     return embeddings
61  |
62  | def generate_market_analysis(llm: LLMWrapper, text: str, source: str) -> str:
63  |     prompt = format_analysis_prompt(
64  |         sector=DEFAULT_SECTOR,
65  |         data=f"Source: {source}\n\nContent: {text}",
66  |         key_metrics=DEFAULT_KEY_METRICS,
67  |     )
68  |
69  |     return llm.query(
70  |         prompt,
71  |         system_prompt=(
72  |             "You are a senior market research analyst. "
73  |             "Extract startup pain points, unmet needs, and market signals."
74  |         ),
75  |     )
76  |
77  | def build_processed_records(raw_records: List[dict]) -> List[dict]:
78  |
79  |     cleaned_texts = [
80  |         preprocess_text(record.get("title"), record.get("body"))
81  |         for record in raw_records
82  |     ]
83  |
84  |     embeddings = generate_embeddings(cleaned_texts)
85  |
86  |     analysis_model = ModelRouter.get_model_for_task("synthesis")
87  |     llm = LLMWrapper(model=analysis_model)
88  |
89  |     processed_records: List[dict] = []
90  |
91  |     for record, cleaned_text, embedding in zip(
92  |         raw_records, cleaned_texts, embeddings
93  |     ):
94  |         processed_records.append(
95  |             {
96  |                 "title": record.get("title"),
```

```
 97 |                    "body": record.get("body"),
 98 |                    "score": record.get("score"),
 99 |                    "comments": record.get("comments"),
100 |                    "created_at": record.get("created_at"),
101 |                    "source": record.get("source"),
102 |                    "cleaned_text": cleaned_text,
103 |                    "embedding": embedding,
104 |                    "market_analysis": generate_market_analysis(
105 |                        llm,
106 |                        cleaned_text,
107 |                        record.get("source", "unknown"),
108 |                    ),
109 |                }
110 |            )
111 |
112 |        return processed_records
113 |
114 |  def run_pipeline() -> Path:
115 |        raw_path = Path(settings.paths.RAW_DATA)
116 |        processed_path = Path(settings.paths.PROCESSED_DATA)
117 |        processed_path.mkdir(parents=True, exist_ok=True)
118 |
119 |        reddit_records = load_json_records(raw_path / "reddit.json")
120 |        hn_records = load_json_records(raw_path / "hn.json")
121 |
122 |        combined_records = reddit_records + hn_records
123 |
124 |        processed_records = build_processed_records(combined_records)
125 |
126 |        output_file = processed_path / "market_analysis.json"
127 |
128 |        with output_file.open("w", encoding="utf-8") as handle:
129 |            json.dump(processed_records, handle, ensure_ascii=False, indent=2)
130 |
131 |        print(f"Processed {len(processed_records)} records -> {output_file}")
132 |
133 |        return output_file
134 |
135 |
136 |  if __name__ == "__main__":
137 |        run_pipeline()
138 |
```

# inference_pipeline.py

```
1 |
2 |
```

# training_pipeline.py

```
1 |
2 |
```

```python
1  | """
2  | Script to generate market intelligence reports
3  |
4  | This script can generate reports in multiple ways:
5  | 1. From processed data (market_analysis.json)
6  | 2. Using the multi-agent system for comprehensive analysis
7  | 3. Using the RAG insight generator
8  |
9  | Usage:
10 |     python scripts/generate_report.py [--method METHOD] [--topic TOPIC] [--output OUTPUT]
11 |
12 | Options:
13 |     --method: 'data' (from processed data), 'agents' (multi-agent), 'rag' (RAG insights)
14 |     --topic: Topic for agent-based or RAG-based reports
15 |     --output: Output file path (default: reports/report_YYYYMMDD_HHMMSS.json)
16 | """
17 |
18 | import os
19 | import sys
20 | import json
21 | import argparse
22 | from pathlib import Path
23 | from datetime import datetime
24 |
25 | # Add project root to path
26 | project_root = Path(__file__).parent.parent
27 | sys.path.insert(0, str(project_root))
28 |
29 | from app.config.settings import get_settings
30 | from app.rag.insight_generator import InsightGenerator
31 | from app.rag.retriever import RAGRetriever
32 | from app.vector_db.config import VectorDBConfig
33 | from app.agents import Orchestrator
34 |
35 | settings = get_settings()
36 |
37 |
38 | def generate_report_from_data(output_path: Path) -> dict:
39 |     """Generate report from processed data"""
40 |     processed_file = Path(settings.paths.PROCESSED_DATA) / "market_analysis.json"
41 |
42 |     if not processed_file.exists():
43 |         raise FileNotFoundError(f"Processed data not found: {processed_file}")
44 |
45 |     with open(processed_file, "r", encoding="utf-8") as f:
46 |         data = json.load(f)
47 |
48 |     # Aggregate insights
49 |     all_analyses = [record.get("market_analysis", "") for record in data]
50 |
51 |     # Count sources
52 |     sources = {}
53 |     for record in data:
54 |         source = record.get("source", "unknown")
55 |         sources[source] = sources.get(source, 0) + 1
56 |
57 |     # Generate summary using LLM
58 |     from app.llm.client import LLMWrapper
59 |     from app.llm.model_router import ModelRouter
60 |
61 |     model = ModelRouter.get_model_for_task("synthesis")
62 |     llm = LLMWrapper(model=model)
63 |
64 |     summary_prompt = f"""Based on {len(data)} market analysis records, generate a comprehensive market intelligence report.
65 |
66 | The report should include:
67 | 1. Executive Summary
68 | 2. Key Market Trends
69 | 3. Pain Points Identified
70 | 4. Opportunities
71 | 5. Recommendations
72 |
73 | Analysis excerpts:
74 | {chr(10).join(all_analyses[:10])}
75 |
76 | Generate a structured report in JSON format:
77 | {{
78 |     "executive_summary": "...",
79 |     "key_trends": ["...", "..."],
80 |     "pain_points": ["...", "..."],
81 |     "opportunities": ["...", "..."],
82 |     "recommendations": ["...", "..."]
83 | }}"""
84 |
85 |     response = llm.query(
86 |         summary_prompt,
87 |         system_prompt="You are a senior market intelligence analyst creating executive reports."
88 |     )
89 |
90 |     try:
91 |         report_data = json.loads(response)
92 |     except json.JSONDecodeError:
93 |         report_data = {
94 |             "executive_summary": response[:500],
95 |             "key_trends": [],
96 |             "pain_points": [],
```

```python
 97 |                    "opportunities": [],
 98 |                    "recommendations": []
 99 |                }
100 |
101 |        # Add metadata
102 |        report = {
103 |            "report_type": "data_analysis",
104 |            "generated_at": datetime.now().isoformat(),
105 |            "data_sources": sources,
106 |            "total_records": len(data),
107 |            "report": report_data
108 |        }
109 |
110 |        return report
111 |
112 |
113 | def generate_report_from_agents(topic: str, output_path: Path) -> dict:
114 |     """Generate report using multi-agent system"""
115 |     print(f"■ Using multi-agent system for topic: {topic}")
116 |
117 |     vector_config = VectorDBConfig(
118 |         persist_directory=settings.vectordb.PERSIST_DIR,
119 |         collection_name=settings.vectordb.COLLECTION_NAME,
120 |         embedding_model=settings.model.EMBEDDING_MODEL
121 |     )
122 |
123 |     orchestrator = Orchestrator(vector_db_config=vector_config)
124 |
125 |     result = orchestrator.run_full_pipeline(
126 |         topic=topic,
127 |         generate_ideas=True,
128 |         evaluate_ideas=True,
129 |         num_ideas=5,
130 |         top_k_research=10
131 |     )
132 |
133 |     # Convert to report format
134 |     report = {
135 |         "report_type": "multi_agent",
136 |         "generated_at": datetime.now().isoformat(),
137 |         "topic": topic,
138 |         "execution_time": result.execution_time,
139 |         "research": {
140 |             "summary": result.research_results.summary if result.research_results else "",
141 |             "key_findings": result.research_results.key_findings if result.research_results else [],
142 |             "sources": result.research_results.sources if result.research_results else []
143 |         } if result.research_results else None,
144 |         "market_analysis": {
145 |             "market_size": result.market_analysis.market_size if result.market_analysis else None,
146 |             "market_maturity": result.market_analysis.market_maturity if result.market_analysis else None,
147 |             "trends": result.market_analysis.trends if result.market_analysis else [],
148 |             "opportunities": result.market_analysis.opportunities if result.market_analysis else [],
149 |             "threats": result.market_analysis.threats if result.market_analysis else []
150 |         } if result.market_analysis else None,
151 |         "generated_ideas": [
152 |             {
153 |                 "title": idea.title,
154 |                 "description": idea.description,
155 |                 "target_audience": idea.target_audience,
156 |                 "value_proposition": idea.value_proposition,
157 |                 "key_features": idea.key_features,
158 |                 "market_opportunity": idea.market_opportunity
159 |             }
160 |             for idea in (result.generated_ideas or [])
161 |         ],
162 |         "evaluations": [
163 |             {
164 |                 "idea": eval_result.idea[:100] + "..." if len(eval_result.idea) > 100 else eval_result.idea,
165 |                 "overall_score": eval_result.overall_score,
166 |                 "feasibility_score": eval_result.feasibility_score,
167 |                 "market_potential_score": eval_result.market_potential_score,
168 |                 "innovation_score": eval_result.innovation_score,
169 |                 "verdict": eval_result.verdict,
170 |                 "strengths": eval_result.strengths,
171 |                 "weaknesses": eval_result.weaknesses,
172 |                 "recommendations": eval_result.recommendations
173 |             }
174 |             for eval_result in (result.evaluations or [])
175 |         ],
176 |         "top_ideas": [
177 |             {
178 |                 "title": idea.title,
179 |                 "description": idea.description,
180 |                 "value_proposition": idea.value_proposition
181 |             }
182 |             for idea in (result.top_ideas or [])
183 |         ],
184 |         "summary": result.summary
185 |     }
186 |
187 |     return report
188 |
189 |
190 | def generate_report_from_rag(topic: str, output_path: Path) -> dict:
191 |     """Generate report using RAG insight generator"""
192 |     print(f"■ Using RAG system for topic: {topic}")
193 |
194 |     vector_config = VectorDBConfig(
```

```
195 |             persist_directory=settings.vectordb.PERSIST_DIR,
196 |             collection_name=settings.vectordb.COLLECTION_NAME,
197 |             embedding_model=settings.model.EMBEDDING_MODEL
198 |         )
199 |
200 |         retriever = RAGRetriever(vector_config)
201 |         insight_generator = InsightGenerator()
202 |
203 |         # Retrieve relevant documents
204 |         documents = retriever.retrieve(topic, top_k=10)
205 |
206 |         # Format context
207 |         context = "\n\n".join([
208 |             f"Document {i+1} (Score: {doc.score:.3f}):\n{doc.text[:500]}..."
209 |             for i, doc in enumerate(documents)
210 |         ])
211 |
212 |         # Generate insights
213 |         insights = insight_generator.generate(topic, context)
214 |
215 |         # Create report
216 |         report = {
217 |             "report_type": "rag_insights",
218 |             "generated_at": datetime.now().isoformat(),
219 |             "topic": topic,
220 |             "documents_analyzed": len(documents),
221 |             "insights": {
222 |                 "summary": insights.summary,
223 |                 "pain_points": insights.pain_points,
224 |                 "opportunities": insights.opportunities,
225 |                 "signals": insights.signals
226 |             },
227 |             "sources": list(set([doc.metadata.get("source", "unknown") for doc in documents]))
228 |         }
229 |
230 |         return report
231 |
232 |
233 | def save_report(report: dict, output_path: Path):
234 |     """Save report to file"""
235 |     output_path.parent.mkdir(parents=True, exist_ok=True)
236 |
237 |     with open(output_path, "w", encoding="utf-8") as f:
238 |         json.dump(report, f, ensure_ascii=False, indent=2)
239 |
240 |     print(f"\n■ Report saved to: {output_path}")
241 |
242 |
243 | def print_report_summary(report: dict):
244 |     """Print a summary of the report"""
245 |     print("\n" + "=" * 60)
246 |     print("Report Summary")
247 |     print("=" * 60)
248 |
249 |     report_type = report.get("report_type", "unknown")
250 |     print(f"\n■ Report Type: {report_type}")
251 |     print(f"■ Generated: {report.get('generated_at', 'Unknown')}")
252 |
253 |     if report_type == "data_analysis":
254 |         print(f"■ Records Analyzed: {report.get('total_records', 0)}")
255 |         report_data = report.get("report", {})
256 |         print(f"■ Trends Identified: {len(report_data.get('key_trends', []))}")
257 |         print(f"■ Opportunities: {len(report_data.get('opportunities', []))}")
258 |
259 |     elif report_type == "multi_agent":
260 |         print(f"■■  Execution Time: {report.get('execution_time', 0):.2f}s")
261 |         if report.get("research"):
262 |             print(f"■ Research Findings: {len(report['research'].get('key_findings', []))}")
263 |         if report.get("market_analysis"):
264 |             print(f"■ Market Maturity: {report['market_analysis'].get('market_maturity', 'Unknown')}")
265 |         print(f"■ Ideas Generated: {len(report.get('generated_ideas', []))}")
266 |         print(f"■ Top Ideas: {len(report.get('top_ideas', []))}")
267 |
268 |     elif report_type == "rag_insights":
269 |         print(f"■ Documents Analyzed: {report.get('documents_analyzed', 0)}")
270 |         insights = report.get("insights", {})
271 |         print(f"■ Pain Points: {len(insights.get('pain_points', []))}")
272 |         print(f"■ Opportunities: {len(insights.get('opportunities', []))}")
273 |         print(f"■ Signals: {len(insights.get('signals', []))}")
274 |
275 |
276 | def main():
277 |     """Main function"""
278 |     parser = argparse.ArgumentParser(description="Generate market intelligence reports")
279 |     parser.add_argument(
280 |         "--method",
281 |         choices=["data", "agents", "rag"],
282 |         default="data",
283 |         help="Report generation method (default: data)"
284 |     )
285 |     parser.add_argument(
286 |         "--topic",
287 |         type=str,
288 |         help="Topic for agent-based or RAG-based reports (required for agents/rag methods)"
289 |     )
290 |     parser.add_argument(
291 |         "--output",
292 |         type=str,
```

```
293 |            help="Output file path (default: reports/report_YYYYMMDD_HHMMSS.json)"
294 |        )
295 |
296 |    args = parser.parse_args()
297 |
298 |    # Validate arguments
299 |    if args.method in ["agents", "rag"] and not args.topic:
300 |        print(f"■ Error: --topic is required for method '{args.method}'")
301 |        return 1
302 |
303 |    # Determine output path
304 |    if args.output:
305 |        output_path = Path(args.output)
306 |    else:
307 |        timestamp = datetime.now().strftime("%Y%m%d_%H%M%S")
308 |        reports_dir = Path(settings.paths.REPORTS)
309 |        output_path = reports_dir / f"report_{timestamp}.json"
310 |
311 |    print("=" * 60)
312 |    print("Market Intelligence Report Generator")
313 |    print("=" * 60)
314 |
315 |    try:
316 |        # Generate report based on method
317 |        if args.method == "data":
318 |            print("\n■ Generating report from processed data...")
319 |            report = generate_report_from_data(output_path)
320 |
321 |        elif args.method == "agents":
322 |            report = generate_report_from_agents(args.topic, output_path)
323 |
324 |        elif args.method == "rag":
325 |            report = generate_report_from_rag(args.topic, output_path)
326 |
327 |        # Save report
328 |        save_report(report, output_path)
329 |
330 |        # Print summary
331 |        print_report_summary(report)
332 |
333 |        return 0
334 |
335 |    except FileNotFoundError as e:
336 |        print(f"\n■ Error: {str(e)}")
337 |        if args.method == "data":
338 |            print("   Run 'python scripts/run_analysis.py' first to process data.")
339 |        return 1
340 |    except Exception as e:
341 |        print(f"\n■ Error generating report: {str(e)}")
342 |        import traceback
343 |        traceback.print_exc()
344 |        return 1
345 |
346 |
347 | if __name__ == "__main__":
348 |    exit(main())
349 |
```

# run_agents.py

```python
"""
Example script demonstrating the Multi-Agent Market Intelligence System

Usage:
    python scripts/run_agents.py
"""

from app.agents import Orchestrator
from app.vector_db.config import VectorDBConfig
from app.config.settings import get_settings

settings = get_settings()


def main():
    """Run the multi-agent system"""

    # Initialize vector DB config
    vector_config = VectorDBConfig(
        persist_directory=settings.vectordb.PERSIST_DIR,
        collection_name=settings.vectordb.COLLECTION_NAME,
        embedding_model=settings.model.EMBEDDING_MODEL
    )

    # Create orchestrator
    orchestrator = Orchestrator(vector_db_config=vector_config)

    # Example: Research and generate ideas for a topic
    topic = "AI-powered productivity tools for remote teams"

    print("=" * 60)
    print("Multi-Agent Market Intelligence System")
    print("=" * 60)
    print(f"\nTopic: {topic}\n")

    # Run full pipeline
    result = orchestrator.run_full_pipeline(
        topic=topic,
        generate_ideas=True,
        evaluate_ideas=True,
        num_ideas=5,
        top_k_research=10
    )

    # Display results
    print("\n" + "=" * 60)
    print("RESULTS SUMMARY")
    print("=" * 60)

    if result.research_results:
        print(f"\n■ Research Results:")
        print(f"   Documents found: {len(result.research_results.documents)}")
        print(f"   Key findings: {len(result.research_results.key_findings)}")
        print(f"   Sources: {', '.join(result.research_results.sources[:3])}...")

    if result.market_analysis:
        print(f"\n■ Market Analysis:")
        print(f"   Market Maturity: {result.market_analysis.market_maturity}")
        print(f"   Market Size: {result.market_analysis.market_size}")
        print(f"   Opportunities: {len(result.market_analysis.opportunities)}")
        print(f"   Trends: {', '.join(result.market_analysis.trends[:3])}...")

    if result.generated_ideas:
        print(f"\n■ Generated Ideas ({len(result.generated_ideas)}):")
        for i, idea in enumerate(result.generated_ideas, 1):
            print(f"\n   {i}. {idea.title}")
            print(f"      {idea.description[:100]}...")
            print(f"      Target: {idea.target_audience}")

    if result.evaluations:
        print(f"\n■ Evaluations:")
        for i, eval_result in enumerate(result.evaluations, 1):
            print(f"   Idea {i}:")
            print(f"      Overall Score: {eval_result.overall_score:.2f}/1.0")
            print(f"      Verdict: {eval_result.verdict}")
            print(f"      Strengths: {len(eval_result.strengths)}")

    if result.top_ideas:
        print(f"\n■ Top {len(result.top_ideas)} Ideas:")
        for i, idea in enumerate(result.top_ideas, 1):
            eval_result = result.evaluations[i-1] if i <= len(result.evaluations) else None
            score = eval_result.overall_score if eval_result else 0.0
            print(f"\n   {i}. {idea.title} (Score: {score:.2f})")
            print(f"      {idea.value_proposition}")

    print(f"\n■■  Execution Time: {result.execution_time:.2f}s")
    print(f"■ Summary: {result.summary}")
    print("\n" + "=" * 60)


if __name__ == "__main__":
    main()
```

```python
"""
Script to run data analysis pipeline

This script:
1. Loads raw scraped data (Reddit and Hacker News)
2. Preprocesses and cleans the text
3. Generates embeddings
4. Runs LLM analysis to extract market insights
5. Saves processed data to data/processed/

Usage:
    python scripts/run_analysis.py
"""

import os
import sys
from pathlib import Path

# Add project root to path
project_root = Path(__file__).parent.parent
sys.path.insert(0, str(project_root))

from pipelines.data_pipeline import run_pipeline, load_json_records
from app.config.settings import get_settings

settings = get_settings()


def main():
    """Run the data analysis pipeline"""
    print("=" * 60)
    print("Data Analysis Pipeline")
    print("=" * 60)

    raw_path = Path(settings.paths.RAW_DATA)
    processed_path = Path(settings.paths.PROCESSED_DATA)

    # Check if raw data exists
    reddit_file = raw_path / "reddit.json"
    hn_file = raw_path / "hn.json"

    files_exist = []
    if reddit_file.exists():
        files_exist.append(("Reddit", reddit_file))
    if hn_file.exists():
        files_exist.append(("Hacker News", hn_file))

    if not files_exist:
        print(f"\n■ No raw data files found in {raw_path}")
        print("   Please run 'python scripts/run_scrapers.py' first to collect data.")
        return 1

    print(f"\n■ Found {len(files_exist)} data source(s):")
    for name, file_path in files_exist:
        try:
            records = load_json_records(file_path)
            print(f"   ■ {name}: {len(records)} records")
        except Exception as e:
            print(f"   ■ {name}: Error loading - {str(e)}")
            return 1

    # Check for OpenAI API key
    if not settings.OPENAI_API_KEY:
        print("\n■ OPENAI_API_KEY not found in settings")
        print("   Please set it in your .env file or environment variables.")
        return 1

    print(f"\n■ Configuration:")
    print(f"   Embedding Model: {settings.model.EMBEDDING_MODEL}")
    print(f"   Analysis Model: {settings.model_routing.ANALYSIS_MODEL}")
    print(f"   Temperature: {settings.model.TEMPERATURE}")
    print(f"   Max Tokens: {settings.model.MAX_TOKENS}")

    # Run pipeline
    try:
        print("\n■ Starting analysis pipeline...")
        print("   This may take a while depending on the amount of data...")

        output_file = run_pipeline()

        print("\n" + "=" * 60)
        print("Pipeline Completed Successfully!")
        print("=" * 60)
        print(f"\n■ Processed data saved to: {output_file}")

        # Show statistics
        try:
            import json
            with open(output_file, "r", encoding="utf-8") as f:
                processed_data = json.load(f)

            print(f"\n■ Statistics:")
            print(f"   Total records processed: {len(processed_data)}")

            # Count by source
            sources = {}
```

```
 97 |                 for record in processed_data:
 98 |                     source = record.get("source", "unknown")
 99 |                     sources[source] = sources.get(source, 0) + 1
100 |
101 |                 print(f"   Sources:")
102 |                 for source, count in sources.items():
103 |                     print(f"      - {source}: {count} records")
104 |
105 |                 print(f"\n■ Next steps:")
106 |                 print(f"   1. Index data into vector DB (if not already done)")
107 |                 print(f"   2. Run 'python scripts/generate_report.py' to generate insights")
108 |                 print(f"   3. Or use the multi-agent system: 'python scripts/run_agents.py'")
109 |
110 |             except Exception as e:
111 |                 print(f"   ■■  Could not load statistics: {str(e)}")
112 |
113 |             return 0
114 |
115 |     except FileNotFoundError as e:
116 |         print(f"\n■ File not found: {str(e)}")
117 |         print("   Make sure raw data files exist in data/raw/")
118 |         return 1
119 |     except Exception as e:
120 |         print(f"\n■ Pipeline failed: {str(e)}")
121 |         import traceback
122 |         traceback.print_exc()
123 |         return 1
124 |
125 |
126 | if __name__ == "__main__":
127 |     exit(main())
128 |
```

# run_scrapers.py

```python
"""
Script to run data collection scrapers (Reddit and Hacker News)

Usage:
    python scripts/run_scrapers.py
"""

import os
import sys
from pathlib import Path

# Add project root to path
project_root = Path(__file__).parent.parent
sys.path.insert(0, str(project_root))

from app.data_collectors.reddit_scraper import scrape_reddit
from app.data_collectors.hn_scraper import scrape_hackernews
from app.config.settings import get_settings

settings = get_settings()


def main():
    """Run all scrapers"""
    print("=" * 60)
    print("Data Collection Scrapers")
    print("=" * 60)

    # Check for required environment variables for Reddit
    reddit_required = ["REDDIT_CLIENT_ID", "REDDIT_CLIENT_SECRET", "REDDIT_USER_AGENT"]
    reddit_missing = [var for var in reddit_required if not os.getenv(var)]

    if reddit_missing:
        print(f"\n■■  Warning: Missing Reddit environment variables: {', '.join(reddit_missing)}")
        print("   Reddit scraping will be skipped.")
        print("   Set these in your .env file or environment.")
        reddit_enabled = False
    else:
        reddit_enabled = True

    results = {}

    # Scrape Reddit
    if reddit_enabled:
        try:
            print("\n■ Scraping Reddit...")
            reddit_path = scrape_reddit(limit=settings.scraper.REDDIT_LIMIT)
            results["reddit"] = {
                "status": "success",
                "path": reddit_path
            }
            print(f"■ Reddit scraping completed")
        except Exception as e:
            print(f"■ Reddit scraping failed: {str(e)}")
            results["reddit"] = {
                "status": "failed",
                "error": str(e)
            }
    else:
        results["reddit"] = {
            "status": "skipped",
            "reason": "Missing environment variables"
        }

    # Scrape Hacker News
    try:
        print("\n■ Scraping Hacker News...")
        hn_path = scrape_hackernews(limit=settings.scraper.HN_LIMIT)
        results["hackernews"] = {
            "status": "success",
            "path": hn_path
        }
        print(f"■ Hacker News scraping completed")
    except Exception as e:
        print(f"■ Hacker News scraping failed: {str(e)}")
        results["hackernews"] = {
            "status": "failed",
            "error": str(e)
        }

    # Summary
    print("\n" + "=" * 60)
    print("Scraping Summary")
    print("=" * 60)

    for source, result in results.items():
        status_icon = "■" if result["status"] == "success" else "■" if result["status"] == "failed" else "■■"
        print(f"{status_icon} {source.capitalize()}: {result['status']}")
        if result["status"] == "success":
            print(f"   Path: {result['path']}")
        elif result["status"] == "failed":
            print(f"   Error: {result.get('error', 'Unknown error')}")

    # Check if we have data for analysis
    success_count = sum(1 for r in results.values() if r["status"] == "success")
    if success_count > 0:
```

```
 97 |            print(f"\n■ Successfully scraped {success_count} source(s)")
 98 |            print("  Next step: Run 'python scripts/run_analysis.py' to process the data")
 99 |        else:
100 |            print("\n■■  No data was successfully scraped. Please check errors above.")
101 |            return 1
102 |
103 |        return 0
104 |
105 |
106 | if __name__ == "__main__":
107 |     exit(main())
108 |
```