# Relevant Code

## Code Repository

**Repository**: Not Available (Official implementation not found in paper)

**Status**: Pseudocode (No Repository Found)

**Language**: Python (Conceptual Implementation)

**Last Updated**: November 2025 (Based on paper publication)

**Stars/Popularity**: N/A

## Architecture Overview

The Kimi K2 architecture consists of three main computational systems: the MuonClip optimizer for stable training of trillion-scale models, the agentic data synthesis pipeline for generating large-scale tool-use demonstrations, and the unified reinforcement learning framework combining verifiable rewards with self-critique mechanisms. The conceptual architecture follows a modular design where the Mixture-of-Experts (MoE) model serves as the core computational engine, supported by specialized subsystems for optimization, data generation, and alignment training.

The system is designed around the principle of token efficiency - maximizing learning signal per training token through optimized optimization algorithms and synthetic data augmentation. Key architectural components include the QK-Clip mechanism for attention stabilization, the multi-head latent attention (MLA) system for efficient computation, and the hierarchical expert routing system that activates 8 out of 384 experts per forward pass.

### Directory Structure

```
kimi_k2_architecture/
├── optimizers/          - MuonClip optimizer implementation
├── models/              - MoE transformer architecture with MLA
├── data_synthesis/      - Agentic data generation pipeline
├── rl_framework/        - Unified RL with verifiable rewards
```

```
├── training/          - Distributed training infrastructure
└── evaluation/        - Comprehensive benchmark evaluation
```

# Key Implementation

## MuonClip Optimizer with QK-Clip

This implements the core innovation from Algorithm 1 in the paper - a stabilized version of the Muon optimizer that prevents attention logit explosion through per-head weight clipping.

```python
class MuonClipOptimizer:
    def __init__(self, params, lr=2e-4, weight_decay=0.1, qk_clip_threshold=100.0):
        self.lr = lr
        self.weight_decay = weight_decay
        self.qk_threshold = qk_clip_threshold
        self.momentum_buffers = {}


    def step(self, model, gradients):
        # Step 1: Standard Muon optimizer update
        for name, param in model.named_parameters():
            if param.grad is None:
                continue

            # Momentum accumulation
            if name not in self.momentum_buffers:
                self.momentum_buffers[name] = torch.zeros_like(param.grad)
            self.momentum_buffers[name] = 0.9 * self.momentum_buffers[name] + param.grad

            # Newton-Schulz iteration for inverse square root
            momentum = self.momentum_buffers[name]
            n, m = momentum.shape
            scale = torch.sqrt(max(n, m)) * 0.2
            preconditioned = newton_schulz(momentum) * scale

            # Weight update with decay
            param.data = param.data - self.lr * (preconditioned + self.weight_decay * param.data)

        # Step 2: QK-Clip for attention stability
        self._apply_qk_clip(model)
```

```python
    def _apply_qk_clip(self, model):
        """Apply per-head QK-Clip to prevent attention logit explosion"""
        for layer in model.transformer.layers:
            attention = layer.attention

            # Calculate max logits for each head from forward pass cache
            max_logits = attention.get_cached_max_logits()  # S_max^h

            # Apply per-head clipping where needed
            for head_idx, logit_max in enumerate(max_logits):
                if logit_max > self.qk_threshold:
                    gamma = self.qk_threshold / logit_max

                    # Apply per-head scaling for MLA components
                    # q^C and k^C (head-specific) scaled by sqrt(gamma)
                    attention.q_proj_head[head_idx].data *= torch.sqrt(gamma)
                    attention.k_proj_head[head_idx].data *= torch.sqrt(gamma)

                    # q^R (head-specific rotary) scaled by gamma
                    attention.q_rotary[head_idx].data *= gamma

                    # k^R (shared rotary) left untouched
```

**Key aspects**:
- Maintains Muon's token efficiency while preventing training instability
- Per-head clipping minimizes intervention on model training
- Selective application to unshared MLA components preserves cross-head interactions
- Threshold $\tau=100$ provides stable training for trillion-scale models

## Agentic Data Synthesis Pipeline

This implements the three-stage data generation system described in Section 3.1.1, creating diverse tool-use demonstrations through simulated environments.

```python
class AgenticDataSynthesisPipeline:
    def __init__(self):
        self.tool_repository = ToolRepository()
        self.agent_generator = AgentGenerator()
        self.trajectory_generator = TrajectoryGenerator()
        self.quality_evaluator = QualityEvaluator()
```

```python
    def generate_training_data(self, num_trajectories=50000):
        """Generate large-scale agentic training data"""
        training_data = []

        for _ in range(num_trajectories):
            # Stage 1: Tool spec generation
            tool_set = self._sample_tool_set()

            # Stage 2: Agent and task generation
            agent = self._generate_agent(tool_set)
            task, rubric = self._generate_task_with_rubric(agent, tool_set)

            # Stage 3: Trajectory generation
            trajectory = self._generate_trajectory(agent, task, tool_set)

            # Quality evaluation and filtering
            if self.quality_evaluator.evaluate(trajectory, rubric):
                training_data.append(trajectory)

        return training_data

    def _sample_tool_set(self):
        """Sample diverse combination of real and synthetic tools"""
        # 50% real MCP tools, 50% synthetic tools
        real_tools = random.sample(self.tool_repository.real_tools, k=2)
        synthetic_tools = random.sample(self.tool_repository.synthetic_tools, k=4)
        return real_tools + synthetic_tools

    def _generate_trajectory(self, agent, task, tool_set):
        """Generate multi-turn tool-use trajectory with simulation"""
        environment = ToolSimulator(tool_set)
        user_simulator = UserSimulator()

        trajectory = Trajectory()
        conversation = user_simulator.initiate_conversation(task)

        max_turns = 10
        for turn in range(max_turns):
            # Agent reasoning and tool selection
            agent_action = agent.reason_and_act(conversation, tool_set)
```

```
            if agent_action.type == "tool_call":
                # Execute tool in simulated environment
                tool_result = environment.execute_tool(
                    agent_action.tool_name,
                    agent_action.parameters
                )

                # Add realistic stochasticity and potential failures
                tool_result = self._add_realistic_variation(tool_result)

                trajectory.add_tool_call(agent_action, tool_result)
                conversation.extend([agent_action, tool_result])

            elif agent_action.type == "final_response":
                trajectory.add_final_response(agent_action.content)
                break

        return trajectory
```

**Key aspects**:
- Combines 3000+ real MCP tools with 20,000+ synthetic tools
- Hierarchical domain evolution ensures comprehensive tool coverage
- Multi-turn simulation with realistic environment feedback
- Quality filtering based on task success criteria and rubric evaluation

## Unified Reinforcement Learning Framework

This implements the dual reward system combining verifiable rewards for structured tasks with self-critique mechanisms for subjective domains.

```
class UnifiedRLFramework:
    def __init__(self, actor_model, critic_model):
        self.actor = actor_model
        self.critic = critic_model
        self.verifiable_envs = VerifiableEnvironments()
        self.rubric_evaluator = RubricEvaluator()

    def train_step(self, batch_size=8, temperature=1.0):
        """Unified RL training step with both reward types"""
        total_loss = 0
```

```python
        for problems in self._sample_batch(batch_size):
            # Generate multiple responses per problem
            responses = []
            for _ in range(batch_size):
                response = self.actor.generate(
                    problems.prompt,
                    temperature=temperature,
                    max_tokens=self._get_token_budget(problems.type)
                )
                responses.append(response)

            # Calculate rewards based on problem type
            if problems.is_verifiable:
                rewards = self._verifiable_rewards(problems, responses)
            else:
                rewards = self._self_critique_rewards(problems, responses)

            # Policy optimization with Muon optimizer
            loss = self._compute_policy_loss(responses, rewards)
            total_loss += loss

            # Update critic model on-policy
            if problems.is_verifiable:
                self._update_critic_verifiable(problems, responses, rewards)

        # Apply temperature decay schedule
        self._update_temperature()

        return total_loss / batch_size

    def _verifiable_rewards(self, problems, responses):
        """Calculate objective rewards for verifiable tasks"""
        rewards = []
        for response in responses:
            if problems.type == "math":
                reward = self.verifiable_envs.math_solver.verify(response, problems.solution)
            elif problems.type == "coding":
                reward = self.verifiable_envs.code_executor.verify(response, problems.test_suite)
            elif problems.type == "instruction_following":
                reward = self.verifiable_envs.constraint_checker.verify(response, problems.constrai
```

```python
                else:
                    reward = self.verifiable_envs.stem_evaluator.verify(response, problems.answer_key)
                rewards.append(reward)
        return rewards


    def _self_critique_rewards(self, problems, responses):
        """Calculate subjective rewards using self-critique"""
        rewards = []

        # Generate pairwise comparisons
        for i, response_i in enumerate(responses):
            response_rewards = []
            for j, response_j in enumerate(responses):
                if i != j:
                    # Use critic model for pairwise comparison
                    comparison = self.critic.compare_responses(
                        problems.prompt,
                        response_i,
                        response_j,
                        rubric=self._get_combined_rubric(problems)
                    )
                    response_rewards.append(comparison.preference_score)

            # Aggregate pairwise scores into final reward
            rewards.append(np.mean(response_rewards))

        return rewards


    def _get_combined_rubric(self, problems):
        """Combine core, prescriptive, and context-specific rubrics"""
        return {
            "core_rubrics": self.rubric_evaluator.get_core_values(),  # Fundamental AI values
            "prescriptive_rubrics": self.rubric_evaluator.get_safety_constraints(),  # Anti-hacking
            "context_rubrics": self.rubric_evaluator.get_domain_specific(problems.domain)
        }
```

**Key aspects**:
- Unified framework handles both verifiable and subjective tasks
- Self-critique mechanism extends alignment to open-ended domains
- Budget control prevents excessive token generation
- Temperature decay balances exploration and exploitation

# Mixture-of-Experts Model Architecture

This implements the 1-trillion parameter MoE architecture with Multi-head Latent Attention (MLA) and efficient expert routing.

```python
class KimiK2MoEModel(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.config = config

        # Model dimensions (from Table 2 in paper)
        self.hidden_size = 7168
        self.num_layers = 61
        self.num_total_experts = 384
        self.num_experts_per_token = 8
        self.num_attention_heads = 64  # Reduced from 128 for efficiency
        self.moe_expert_hidden_size = 2048

        # Embeddings and position encodings
        self.embed_tokens = nn.Embedding(config.vocab_size, self.hidden_size)
        self.rotary_emb = RotaryEmbedding(self.hidden_size // self.num_attention_heads)

        # Transformer layers with MoE
        self.layers = nn.ModuleList([
            TransformerLayer(self.config) for _ in range(self.num_layers)
        ])

        # Output projection
        self.lm_head = nn.Linear(self.hidden_size, config.vocab_size, bias=False)

    def forward(self, input_ids, attention_mask=None):
        hidden_states = self.embed_tokens(input_ids)

        # Apply rotary position embeddings
        for layer in self.layers:
            hidden_states = layer(
                hidden_states,
                attention_mask=attention_mask,
                rotary_emb=self.rotary_emb
            )
```

```python
        logits = self.lm_head(hidden_states)
        return logits


class TransformerLayer(nn.Module):
    def __init__(self, config):
        super().__init__()
        # Multi-head Latent Attention (MLA)
        self.attention = MultiHeadLatentAttention(config)

        # MoE Feed-forward with 384 total experts, 8 active per token
        self.moe = MixtureOfExperts(
            total_experts=384,
            experts_per_token=8,
            expert_hidden_size=2048,
            hidden_size=config.hidden_size
        )

        self.attention_norm = nn.RMSNorm(config.hidden_size)
        self.ffn_norm = nn.RMSNorm(config.hidden_size)


class MultiHeadLatentAttention(nn.Module):
    def __init__(self, config):
        super().__init__()
        self.hidden_size = config.hidden_size
        self.num_heads = config.num_attention_heads
        self.head_dim = self.hidden_size // self.num_heads

        # Latent attention projections
        self.q_proj = nn.Linear(self.hidden_size, self.hidden_size, bias=False)
        self.k_proj = nn.Linear(self.hidden_size, self.hidden_size, bias=False)
        self.v_proj = nn.Linear(self.hidden_size, self.hidden_size, bias=False)
        self.o_proj = nn.Linear(self.hidden_size, self.hidden_size, bias=False)

        # Components for QK-Clip (separate for per-head clipping)
        self.q_proj_head = nn.ModuleList([
            nn.Linear(self.hidden_size, self.head_dim, bias=False)
            for _ in range(self.num_heads)
        ])
        self.k_proj_head = nn.ModuleList([
            nn.Linear(self.hidden_size, self.head_dim, bias=False)
            for _ in range(self.num_heads)
```

```python
        ])


class MixtureOfExperts(nn.Module):
    def __init__(self, total_experts, experts_per_token, expert_hidden_size, hidden_size):
        super().__init__()
        self.total_experts = total_experts
        self.experts_per_token = experts_per_token

        # Gating mechanism for expert routing
        self.gate = nn.Linear(hidden_size, total_experts, bias=False)

        # Expert networks (shared expert + routed experts)
        self.experts = nn.ModuleList([
            FeedForwardExpert(expert_hidden_size, hidden_size)
            for _ in range(total_experts + 1)  # +1 for shared expert
        ])

    def forward(self, hidden_states):
        batch_size, seq_len, hidden_dim = hidden_states.shape

        # Compute gating scores
        gate_logits = self.gate(hidden_states)
        gate_probs = F.softmax(gate_logits, dim=-1)

        # Select top-k experts per token
        top_k_probs, top_k_indices = torch.topk(
            gate_probs,
            self.experts_per_token,
            dim=-1
        )

        # Normalize top-k probabilities
        top_k_probs = top_k_probs / top_k_probs.sum(dim=-1, keepdim=True)

        # Expert computation with load balancing
        expert_outputs = self._compute_expert_outputs(
            hidden_states, top_k_probs, top_k_indices
        )

        return expert_outputs
```

**Key aspects**:
- 1.04 trillion total parameters with 32.6 billion activated parameters
- Sparsity ratio of 48:1 (8 active out of 384 experts)
- Multi-head Latent Attention for computational efficiency
- Per-head attention components enable QK-Clip optimization

# Relation to Paper

| Paper Section | Code Component | Notes |
|---|---|---|
| Section 2.1: MuonClip | `optimizers/ MuonClipOptimizer.py:MuonClipOptimizer` | Implements Algorithm 1 with per-head QK-Clip |
| Section 2.3: Model Architecture | `models/K2MoEModel.py:KimiK2MoEModel` | 1T parameter MoE with MLA attention |
| Section 3.1.1: Data Synthesis | `data_synthesis/ AgenticPipeline.py:AgenticDataSynthesisPipeline` | Three-stage tool-use data generation |
| Section 3.2: Unified RL | `rl_framework/UnifiedRL.py:UnifiedRLFramework` | Verifiable rewards + self-critique |
| Equation QK-Clip | `optimizers/MuonClipOptimizer.py:_apply_qk_clip` | Per-head weight rescaling implementation |
| Figure 4: Rephrasing Pipeline | `data_synthesis/RephrasingPipeline.py` | Chunk-wise autoregressive rephrasing |
| Table 2: Architecture | `models/K2MoEModel.py` | 384 experts, 64 heads, MLA attention |

## Key Differences from Paper

- **Implementation Simplification**: The pseudocode abstracts away distributed training details for clarity

- **Memory Management**: Actual implementation would require sophisticated memory optimization for 1T parameter models

- **Infrastructure Dependencies**: Real implementation requires custom CUDA kernels and distributed computing frameworks

- **Evaluation Integration**: Paper references extensive benchmark evaluation that would require separate evaluation infrastructure