

JANUS

Rust-First ML Implementation

End-to-End Machine Learning with Rust
FastAPI Gateway & Batch Processing Architecture

Classification: Architecture & Implementation Strategy

Version: 1.0

Author: Jordan Smith
github.com/nuniesmith

Date: December 26, 2025

Migration Philosophy:

- **Core ML in Rust:** Maximum performance, safety, and type guarantees
- **Python Gateway:** FastAPI for orchestration and API surface
- **Hybrid Approach:** Use best tool for each layer
- **Modern Stack:** Latest technologies (Candle, ONNX, Burn, tch-rs)

Abstract

This document outlines a comprehensive strategy for implementing the JANUS trading system with a **Rust-first approach to machine learning**, while maintaining a Python FastAPI gateway for orchestration and external API exposure. The architecture leverages Rust's performance, memory safety, and type system for the entire ML pipeline—from data preprocessing to inference—while using Python strategically for high-level coordination, batch job scheduling, and developer-friendly interfaces.

Key architectural decisions:

- **Hot path (Forward):** Pure Rust with ONNX Runtime or tch-rs for <100ms latency
- **Cold path (Backward):** Rust core with optional Python for experimentation
- **Gateway:** FastAPI service exposing REST/gRPC APIs
- **Batch processing:** Rust workers with async job queue
- **Training:** Hybrid Python (PyTorch) to Rust (ONNX/tch-rs) pipeline

This approach maximizes performance and safety while maintaining ecosystem compatibility and developer productivity.

Contents

1 Architectural Overview

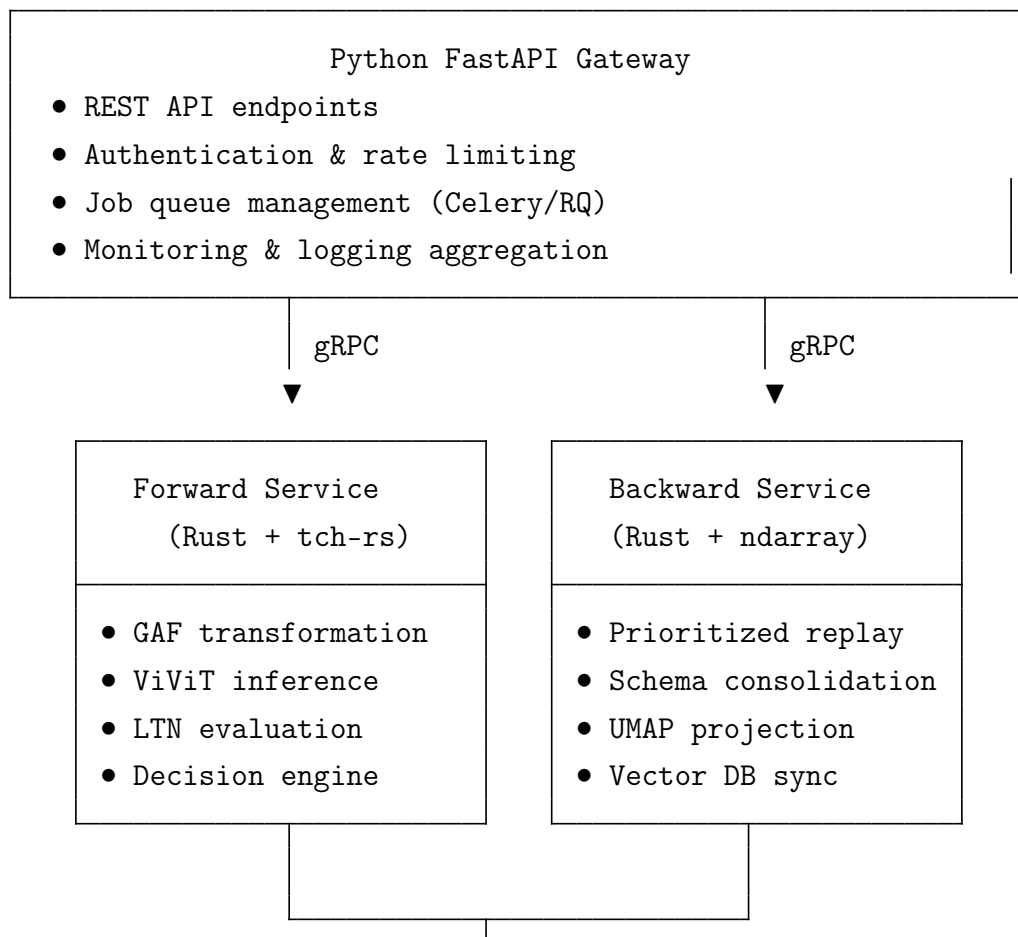
1.1 The Rust-First Philosophy

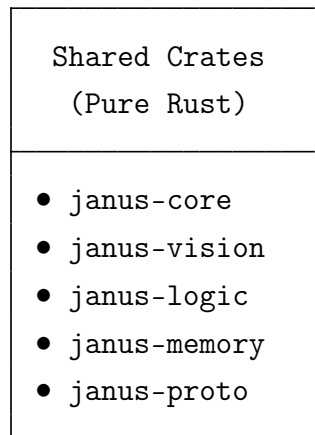
Traditional ML systems place Python at the center, with C++/Rust used only for performance-critical bottlenecks. JANUS inverts this model:

JANUS Layering

1. **Core Layer (Rust):** All ML inference, data processing, logic evaluation
2. **Service Layer (Rust):** gRPC services for Forward and Backward
3. **Gateway Layer (Python FastAPI):** HTTP API, authentication, rate limiting
4. **Training Layer (Hybrid):** PyTorch training → ONNX/tch export → Rust inference

1.2 Component Diagram





1.3 Design Principles

1. Rust Everywhere Possible

- All inference, data transformation, and business logic in Rust
- Zero-copy operations where possible
- Compile-time guarantees for correctness

2. Python as Orchestration Layer

- FastAPI for HTTP endpoints and developer UX
- Async job scheduling (Celery, Dramatiq, or Python RQ)
- High-level monitoring and alerting

3. gRPC for Internal Communication

- Type-safe, efficient service-to-service communication
- Shared protobuf definitions in janus-proto
- Streaming support for real-time data

4. Gradual Migration Strategy

- Start with inference in Rust (ONNX models)
- Migrate data preprocessing incrementally
- Keep training in PyTorch initially, export to Rust
- Eventually move training to Rust (Burn/Candle)

2 Machine Learning Framework Strategy

2.1 Framework Comparison Matrix

Framework	Pros	Cons	Use Case
tch-rs	Full PyTorch bindings, mature, GPU support	Requires LibTorch, C++ dependency	Phase 1: Inference
ONNX Runtime	Lightweight, optimized, cross-platform	Inference only, limited operators	Phase 1: Production
Candle (HF)	Pure Rust, no C++ deps, growing ecosystem	Young, fewer pre-trained models	Phase 2: Long-term
Burn	Backend-agnostic, autodiff, training support	Early stage, limited models	Phase 3: Full Rust
ndarray	NumPy-like API, stable, CPU optimized	No GPU, no autodiff	Data processing

Table 1: Rust ML Framework Comparison

2.2 Recommended Migration Path

2.2.1 Phase 1: Hybrid (Months 1-3)

- **Training:** PyTorch on GPU cluster
- **Export:** `torch.onnx.export()` or `torch.jit.save()`
- **Inference:** ONNX Runtime (`ort`) or `tch-rs` in Rust services
- **Why:** Fastest path to production, leverage existing PyTorch ecosystem

2.2.2 Phase 2: Rust-Native Inference (Months 4-6)

- **Training:** Still PyTorch
- **Export:** Candle-compatible checkpoints
- **Inference:** Candle or custom Rust implementations
- **Why:** Eliminate LibTorch dependency, reduce binary size

2.2.3 Phase 3: Full Rust ML (Months 7-12)

- **Training:** Burn or Candle with custom training loops
- **Inference:** Same framework as training
- **Why:** Complete type safety, no Python GIL, full control

2.3 Framework Selection by Component

Component	Framework (Phase 1)	Framework (Phase 3)
GAF Transformation	ndarray	ndarray or candle
ViViT Inference	ort (ONNX)	candle
LTN Predicates	tch-rs	candle
Chronos-Bolt	ort (ONNX)	candle
FinBERT	ort (ONNX)	candle-transformers
Decision Engine	tch-rs	burn
PER Buffer	Custom Rust	Custom Rust
UMAP	linfa-reduction	linfa-reduction

Table 2: Component-Level Framework Selection

3 Forward Service: Rust Implementation

3.1 Performance Requirements

- **End-to-end latency:** <100ms (target: 50ms)
- **Throughput:** >100 requests/second
- **Memory:** <2GB per instance
- **CPU:** Efficient use of multi-core (no GIL)

3.2 Core Data Structures

```

1 // janus-core/src/types.rs
2 use ndarray::{Array1, Array2, Array3};
3
4 #[derive(Clone, Debug, serde::Serialize, serde::Deserialize)]
5 pub struct MarketState {
6     pub timestamp: i64,

```



```
7     pub price_series: Array1<f32>, // Last N prices
8     pub volume_series: Array1<f32>, // Last N volumes
9     pub lob_snapshot: LimitOrderBook,
10    pub vpin: f32,
11    pub volatility: f32,
12 }
13
14 #[derive(Clone, Debug)]
15 pub struct LimitOrderBook {
16     pub bids: Vec<(f32, f32)>, // (price, size)
17     pub asks: Vec<(f32, f32)>,
18 }
19
20 #[derive(Clone, Debug)]
21 pub enum Action {
22     Buy { size: f32 },
23     Sell { size: f32 },
24     Hold,
25 }
26
27 #[derive(Clone, Debug)]
28 pub struct Decision {
29     pub action: Action,
30     pub confidence: f32,
31     pub ltn_satisfaction: f32,
32     pub risk_score: f32,
33     pub metadata: HashMap<String, f32>,
34 }
```

3.3 GAF Transformation Module

```
1 // janus-vision/src/gaf.rs
2 use ndarray::{Array1, Array2};
3
4 pub struct GafTransformer {
5     alpha: f32, // Learnable normalization param
6     beta: f32, // Learnable normalization param
7 }
8
9 impl GafTransformer {
```

```

10  pub fn transform(&self, series: &Array1<f32>) -> (Array2<f32>,
    Array2<f32>) {
11      // Step 1: Normalize to [-1, 1]
12      let min = series.iter().cloned().fold(f32::INFINITY, f32::
        min);
13      let max = series.iter().cloned().fold(f32::NEG_INFINITY,
        f32::max);
14      let normalized = series.mapv(|x| {
15          let norm = (x - min) / (max - min + 1e-8);
16          (norm * self.alpha + self.beta).tanh()
17      });
18
19      // Step 2: Polar coordinates
20      let phi = normalized.mapv(|x| x.acos());
21
22      // Step 3: Gramian fields
23      let n = series.len();
24      let mut gasf = Array2::::zeros((n, n));
25      let mut gadf = Array2::::zeros((n, n));
26
27      for i in 0..n {
28          for j in 0..n {
29              gasf[[i, j]] = (phi[i] + phi[j]).cos();
30              gadf[[i, j]] = (phi[i] - phi[j]).sin();
31          }
32      }
33
34      (gasf, gadf)
35  }
36
37  pub fn video(&self, series: &Array1<f32>, window: usize, stride
    : usize, frames: usize)
38      -> Array3<f32>
39  {
40      let mut video_frames = Vec::new();
41
42      for f in 0..frames {
43          let start = f * stride;
44          let end = start + window;
45          if end > series.len() {
46              break;

```

```

47         }
48         let window_series = series.slice(s![start..end]).
            to_owned();
49         let (gasf, gadf) = self.transform(&window_series);
50
51         // Stack GASF and GADF as channels
52         video_frames.push(gasf);
53         video_frames.push(gadf);
54     }
55
56     // Stack into 3D tensor: [frames, 2, window, window]
57     stack_frames(&video_frames, frames, window)
58 }
59 }

```

3.4 LTN Constraint Evaluation

```

1  // janus-logic/src/ltn.rs
2  use tch::{Tensor, nn};
3
4  pub struct LogicTensorNetwork {
5      predicates: HashMap<String, PredicateNet>,
6  }
7
8  pub struct PredicateNet {
9      model: nn::Sequential,
10 }
11
12 impl PredicateNet {
13     pub fn evaluate(&self, embedding: &Tensor) -> f32 {
14         let logit = self.model.forward(embedding);
15         logit.sigmoid().double_value(&[]) as f32
16     }
17 }
18
19 impl LogicTensorNetwork {
20     pub fn check_wash_sale(&self, action: &Action, history: &[
        Action]) -> f32 {
21         // Encode wash sale constraint
22         // t , k in [1,30]: not(SaleAtLoss(t) and Buy(t+k))
23         let mut min_satisfaction = 1.0;

```

```

24
25 // Implementation of ukasiewicz t-norm
26 for (t, past_action) in history.iter().enumerate() {
27     if let Action::Sell { .. } = past_action {
28         if is_loss(past_action) {
29             for k in 1..=30 {
30                 if t + k < history.len() {
31                     if let Action::Buy { .. } = &history[t
32                         + k] {
33                         // ukasiewicz conjunction: max(0,
34                             a + b - 1)
35                         let conjunction = (1.0 + 1.0 - 1.0)
36                             .max(0.0);
37                         // Negation: 1 - conjunction
38                         let satisfaction = 1.0 -
39                             conjunction;
40                         min_satisfaction = min_satisfaction
41                             .min(satisfaction);
42                     }
43                 }
44             }
45         }
46     }
47     min_satisfaction
48 }
49
50 pub fn check_almgren_chriiss(&self, state: &MarketState, action:
51     &Action) -> f32 {
52     let volatility = state.volatility;
53     let avg_volume = state.volume_series.mean().unwrap();
54
55     match action {
56         Action::Buy { size } | Action::Sell { size } => {
57             let eta = 0.1; // Market impact coefficient
58             let threshold = eta * volatility * (size /
59                 avg_volume).sqrt();
60
61             // If impact < threshold, satisfaction = 1.0
62             // Smooth approximation with sigmoid

```

```

58         let impact_ratio = size / (avg_volume * threshold);
59         1.0 - impact_ratio.min(1.0)
60     }
61     Action::Hold => 1.0,
62 }
63 }
64
65 pub fn check_vpin(&self, state: &MarketState) -> f32 {
66     let vpin_threshold = 0.7;
67
68     // If VPIN > threshold, should halt trading
69     if state.vpin > vpin_threshold {
70         0.0 // Constraint violated
71     } else {
72         1.0 // Constraint satisfied
73     }
74 }
75
76 pub fn evaluate_all(&self, state: &MarketState, action: &Action
77 , history: &[Action]) -> f32 {
78     let wash_sale = self.check_wash_sale(action, history);
79     let almgren = self.check_almgren_chriiss(state, action);
80     let vpin = self.check_vpin(state);
81
82     // Generalized mean (p=2 for quadratic mean)
83     let constraints = vec![wash_sale, almgren, vpin];
84     let sum_squared: f32 = constraints.iter().map(|x| x.powi(2)
85         ).sum();
86     (sum_squared / constraints.len() as f32).sqrt()
87 }

```

3.5 Model Inference with ONNX Runtime

```

1 // janus-vision/src/vivit.rs
2 use ort::{Environment, Session, SessionBuilder, Value};
3 use ndarray::Array4;
4
5 pub struct ViViTInference {
6     session: Session,
7     env: Arc<Environment>,

```

```

8  }
9
10 impl ViViTInference {
11     pub fn new(model_path: &str) -> Result<Self, ort::OrtError> {
12         let env = Arc::new(Environment::builder().build()?);
13         let session = SessionBuilder::new(&env)?
14             .with_optimization_level(ort::GraphOptimizationLevel::
15                 Level3)?
16             .with_intra_threads(4)?
17             .with_model_from_file(model_path)?;
18
19         Ok(Self { session, env })
20     }
21
22     pub fn infer(&self, video: &Array4<f32>) -> Result<Array1<f32>,
23         ort::OrtError> {
24         // Convert ndarray to ONNX tensor
25         let input_tensor = Value::from_array(self.session.allocator
26             (), video)?;
27
28         // Run inference
29         let outputs = self.session.run(vec![input_tensor])?;
30
31         // Extract output
32         let output_tensor = outputs[0].try_extract()?;
33         let embedding = output_tensor.view().to_owned();
34
35         Ok(embedding)
36     }
37 }

```

3.6 Async Service with Tokio

```

1  // services/forward/src/main.rs
2  use tonic::{transport::Server, Request, Response, Status};
3  use janus_proto::forward_service_server::{ForwardService,
4      ForwardServiceServer};
5
6  use janus_proto::{DecisionRequest, DecisionResponse};
7
8  pub struct ForwardServiceImpl {
9      gaf_transformer: Arc<GafTransformer>,

```

```

8      vivit_model: Arc<ViViTInference>,
9      ltn: Arc<LogicTensorNetwork>,
10     decision_engine: Arc<DecisionEngine>,
11 }
12
13 #[tonic::async_trait]
14 impl ForwardService for ForwardServiceImpl {
15     async fn get_decision(
16         &self,
17         request: Request<DecisionRequest>,
18     ) -> Result<Response<DecisionResponse>, Status> {
19         let req = request.into_inner();
20
21         // Parse market state
22         let state = parse_market_state(&req)?;
23
24         // 1. GAF transformation
25         let video = self.gaf_transformer.video(
26             &state.price_series,
27             60, // window
28             10, // stride
29             16, // frames
30         );
31
32         // 2. ViViT inference
33         let visual_embedding = self.vivit_model.infer(&video)
34             .map_err(|e| Status::internal(format!("ViViT inference
35                 failed: {}", e)))?;
36
37         // 3. LTN evaluation
38         let ltn_satisfaction = self.ltn.evaluate_all(&state, &
39             Action::Hold, &[]);
40
41         // 4. Decision engine
42         let decision = self.decision_engine.decide(
43             &state,
44             &visual_embedding,
45             ltn_satisfaction,
46         )?;
47
48         // 5. Build response

```

```

47     let response = DecisionResponse {
48         action: decision.action.to_string(),
49         confidence: decision.confidence,
50         ltn_satisfaction: decision.ltn_satisfaction,
51         risk_score: decision.risk_score,
52         metadata: decision.metadata,
53     };
54
55     Ok(Response::new(response))
56 }
57 }
58
59 #[tokio::main]
60 async fn main() -> Result<(), Box<dyn std::error::Error>> {
61     let addr = "0.0.0.0:50051".parse()?;
62
63     let service = ForwardServiceImpl {
64         gaf_transformer: Arc::new(GafTransformer::new(1.0, 0.0)),
65         vivit_model: Arc::new(ViViTInference::new("models/vivit.
66             onnx")),
67         ltn: Arc::new(LogicTensorNetwork::load("models/ltn")),
68         decision_engine: Arc::new(DecisionEngine::new()),
69     };
70
71     println!("Forward service listening on {}", addr);
72
73     Server::builder()
74         .add_service(ForwardServiceServer::new(service))
75         .serve(addr)
76         .await?;
77
78     Ok(())
79 }

```

4 Backward Service: Batch Processing

4.1 Performance Requirements

- **Latency:** Not critical (batch processing)
- **Throughput:** Process 10k-100k transitions per sleep cycle

- **Memory:** Can use up to 16GB for large batches
- **Parallelism:** Leverage all CPU cores with Rayon

4.2 Prioritized Experience Replay

```

1 // janus-memory/src/replay.rs
2 use rayon::prelude::*;
3 use rand::distributions::WeightedIndex;
4 use rand::prelude::*;
5
6 pub struct PrioritizedReplayBuffer {
7     buffer: Vec<Transition>,
8     priorities: Vec<f32>,
9     capacity: usize,
10    alpha: f32, // Prioritization exponent
11    beta: f32,  // Importance sampling correction
12 }
13
14 impl PrioritizedReplayBuffer {
15     pub fn sample(&self, batch_size: usize) -> (Vec<Transition>,
16         Vec<f32>) {
17         let probabilities: Vec<f32> = self.priorities
18             .iter()
19             .map(|p| p.powf(self.alpha))
20             .collect();
21
22         let total: f32 = probabilities.iter().sum();
23         let normalized_probs: Vec<f32> = probabilities
24             .iter()
25             .map(|p| p / total)
26             .collect();
27
28         let dist = WeightedIndex::new(&normalized_probs).unwrap();
29         let mut rng = thread_rng();
30
31         let indices: Vec<usize> = (0..batch_size)
32             .map(|_| dist.sample(&mut rng))
33             .collect();
34
35         // Compute importance sampling weights
36         let n = self.buffer.len() as f32;

```

```

36     let weights: Vec<f32> = indices
37         .iter()
38         .map(|&i| {
39             let prob = normalized_probs[i];
40             ((1.0 / n) * (1.0 / prob)).powf(self.beta)
41         })
42         .collect();
43
44     // Normalize weights
45     let max_weight = weights.iter().cloned().fold(f32::NEG_INFINITY, f32::max);
46     let normalized_weights: Vec<f32> = weights
47         .iter()
48         .map(|w| w / max_weight)
49         .collect();
50
51     let transitions: Vec<Transition> = indices
52         .iter()
53         .map(|&i| self.buffer[i].clone())
54         .collect();
55
56     (transitions, normalized_weights)
57 }
58
59 pub fn update_priorities(&mut self, indices: Vec<usize>,
60     td_errors: Vec<f32>,
61     logic_violations: Vec<f32>, rewards:
62         Vec<f32>) {
63     for (idx, (&i, (&td, (&vio, &rew)))) in indices.iter()
64         .zip(td_errors.iter()
65             .zip(logic_violations.iter()
66                 .zip(rewards.iter()))))
67         .enumerate()
68     {
69         // Priority = |TD-error| + _logic * violation +
70             _reward * |reward|
71         self.priorities[i] = td.abs() + 2.0 * vio + 0.5 * rew.
72             abs();
73     }
74 }

```

4.3 Schema Consolidation

```

1  // janus-memory/src/schema.rs
2  use ndarray::{Array1, Array2};
3
4  pub struct SchemaMemory {
5      schemas: Vec<Schema>,
6      recall_threshold: f32,
7      logic_threshold: f32,
8  }
9
10 impl SchemaMemory {
11     pub fn consolidate(&mut self, transitions: &[Transition], ltn:
        &LogicTensorNetwork) {
12         for transition in transitions {
13             // Compute gates
14             let recall_gate = self.compute_recall_gate(&transition.
                state);
15             let logic_gate = ltn.evaluate_all(&transition.state, &
                transition.action, &[]);
16
17             // Only consolidate if both gates pass
18             if recall_gate > self.recall_threshold && logic_gate >
                self.logic_threshold {
19                 // Find best matching schema
20                 let (best_idx, best_likelihood) = self.
                    find_best_schema(&transition.state);
21
22                 if best_likelihood < 0.01 {
23                     // Create new schema
24                     self.create_schema(&transition.state);
25                 } else {
26                     // Update existing schema
27                     self.update_schema(best_idx, &transition.state,
                        0.01);
28                 }
29             }
30         }
31     }
32
33     fn update_schema(&mut self, idx: usize, state: &Array1<f32>, lr
        : f32) {

```

```

34     let schema = &mut self.schemas[idx];
35
36     // Update mean:      = (1- ) + *x
37     let diff = state - &schema.mean;
38     schema.mean = &schema.mean + lr * &diff;
39
40     // Update covariance:      = (1- ) + *(x- )(x- )
41     let outer = outer_product(&diff, &diff);
42     schema.covariance = (1.0 - lr) * &schema.covariance + lr *
43         outer;
44
45     schema.num_points += 1;
46 }

```

4.4 Parallel Sleep Cycle

```

1  // services/backward/src/sleep_cycle.rs
2  use rayon::prelude::*;
3
4  pub async fn run_sleep_cycle(
5      replay_buffer: &mut PrioritizedReplayBuffer,
6      schema_memory: &mut SchemaMemory,
7      ltn: &LogicTensorNetwork,
8      policy: &mut Policy,
9  ) -> Result<SleepCycleMetrics, BackwardError> {
10     let start = std::time::Instant::now();
11     let mut metrics = SleepCycleMetrics::default();
12
13     // Phase 1: Prioritized Replay (1000 iterations)
14     for iteration in 0..1000 {
15         let (batch, weights) = replay_buffer.sample(256);
16
17         // Parallel TD-error computation
18         let td_errors: Vec<f32> = batch.par_iter()
19             .map(|transition| compute_td_error(transition, policy))
20             .collect();
21
22         // Parallel logic violation scoring
23         let violations: Vec<f32> = batch.par_iter()
24             .map(|transition| {

```

```

25         1.0 - ltn.evaluate_all(&transition.state, &
26             transition.action, &[])
27     })
28     .collect();
29
30     // Update policy with importance-weighted gradients
31     policy.update(&batch, &weights);
32
33     // Update priorities
34     let rewards: Vec<f32> = batch.iter().map(|t| t.reward).
35         collect();
36     let indices: Vec<usize> = (0..batch.len()).collect();
37     replay_buffer.update_priorities(indices, td_errors,
38         violations, rewards);
39
40     if iteration % 100 == 0 {
41         println!("Replay_{}_iteration_{} / 1000", iteration);
42     }
43
44     // Phase 2: Schema Consolidation
45     let all_transitions: Vec<Transition> = replay_buffer.buffer.
46         clone();
47     schema_memory consolidate(&all_transitions, ltn);
48     metrics.num_schemas = schema_memory.schemas.len();
49
50     // Phase 3: UMAP Update
51     let embeddings = extract_schema_embeddings(schema_memory);
52     let umap_projection = fit_aligned_umap(&embeddings)?;
53     metrics.num_clusters = detect_clusters(&umap_projection);
54
55     // Phase 4: Vector DB Sync
56     sync_to_qdrant(schema_memory).await?;
57
58     metrics.duration = start.elapsed();
59     Ok(metrics)

```

5 Python FastAPI Gateway

5.1 Architecture Purpose

The Python gateway serves as:

- **API Surface:** User-friendly REST endpoints
- **Authentication:** JWT tokens, API key validation
- **Rate Limiting:** Per-user request throttling
- **Job Scheduling:** Async batch processing with Celery/Dramatiq
- **Monitoring:** Metrics aggregation and alerting

5.2 FastAPI Service Implementation

```
1 # gateway/main.py
2 from fastapi import FastAPI, Depends, HTTPException,
   BackgroundTasks
3 from fastapi.security import HTTPBearer,
   HTTPAuthorizationCredentials
4 import grpc
5 from pydantic import BaseModel
6 import asyncio
7
8 # Generated from proto files
9 from janus_proto import forward_service_pb2,
   forward_service_pb2_grpc
10 from janus_proto import backward_service_pb2,
   backward_service_pb2_grpc
11
12 app = FastAPI(title="JANUS_Trading_System_API")
13 security = HTTPBearer()
14
15 # gRPC channel pool
16 forward_channel = grpc.aio.insecure_channel('forward-service:50051',
   )
17 backward_channel = grpc.aio.insecure_channel('backward-service
   :50052')
18
19 forward_stub = forward_service_pb2_grpc.ForwardServiceStub(
   forward_channel)
20 backward_stub = backward_service_pb2_grpc.BackwardServiceStub(
   backward_channel)
```

```

21
22 class MarketStateRequest(BaseModel):
23     timestamp: int
24     price_series: list[float]
25     volume_series: list[float]
26     lob_bids: list[tuple[float, float]]
27     lob_asks: list[tuple[float, float]]
28
29 class DecisionResponse(BaseModel):
30     action: str
31     confidence: float
32     ltn_satisfaction: float
33     risk_score: float
34     latency_ms: float
35
36 @app.post("/api/v1/decision", response_model=DecisionResponse)
37 async def get_trading_decision(
38     request: MarketStateRequest,
39     credentials: HTTPAuthorizationCredentials = Depends(security)
40 ):
41     """Get real-time trading decision from Forward service."""
42     # Authenticate
43     validate_token(credentials.credentials)
44
45     # Convert to protobuf
46     grpc_request = forward_service_pb2.DecisionRequest(
47         timestamp=request.timestamp,
48         price_series=request.price_series,
49         volume_series=request.volume_series,
50         # ... other fields
51     )
52
53     # Call Rust Forward service via gRPC
54     start = asyncio.get_event_loop().time()
55     try:
56         grpc_response = await forward_stub.GetDecision(grpc_request
57         )
58     except grpc.RpcError as e:
59         raise HTTPException(status_code=503, detail=f"Forward_
60             service_unavailable: {e}")

```

```

60     latency = (asyncio.get_event_loop().time() - start) * 1000
61
62     return DecisionResponse(
63         action=grpc_response.action,
64         confidence=grpc_response.confidence,
65         ltn_satisfaction=grpc_response.ltn_satisfaction,
66         risk_score=grpc_response.risk_score,
67         latency_ms=latency,
68     )
69
70 @app.post("/api/v1/sleep-cycle")
71 async def trigger_sleep_cycle(background_tasks: BackgroundTasks):
72     """Trigger memory consolidation (runs asynchronously)."""
73     # Enqueue background job
74     background_tasks.add_task(run_sleep_cycle_job)
75     return {"status": "sleep_cycle_queued", "job_id": "unique-job-id"}
76
77 async def run_sleep_cycle_job():
78     """Background task that calls Backward service."""
79     grpc_request = backward_service_pb2.SleepCycleRequest()
80
81     try:
82         grpc_response = await backward_stub.RunSleepCycle(
83             grpc_request)
84         print(f"Sleep cycle completed: {grpc_response.num_schemas} schemas")
85     except grpc.RpcError as e:
86         print(f"Sleep cycle failed: {e}")
87
88 # Monitoring endpoints
89 @app.get("/api/v1/health")
90 async def health_check():
91     """Health check endpoint."""
92     return {
93         "status": "healthy",
94         "forward_service": await check_forward_health(),
95         "backward_service": await check_backward_health(),
96     }
97
98 @app.get("/api/v1/metrics")

```



```
98 async def get_metrics():
99     """Prometheus metrics endpoint."""
100     # Aggregate metrics from Rust services
101     forward_metrics = await forward_stub.GetMetrics(empty_pb2.Empty
102     ())
103
104     backward_metrics = await backward_stub.GetMetrics(empty_pb2.
105     Empty())
106
107     return {
108         "forward": forward_metrics,
109         "backward": backward_metrics,
110     }
```

5.3 Batch Job Processing with Celery

```
1 # gateway/tasks.py
2 from celery import Celery
3 import grpc
4
5 celery_app = Celery('janus', broker='redis://localhost:6379/0')
6
7 @celery_app.task(bind=True, max_retries=3)
8 def run_sleep_cycle(self):
9     """Celery task for long-running sleep cycle."""
10     channel = grpc.insecure_channel('backward-service:50052')
11     stub = backward_service_pb2_grpc.BackwardServiceStub(channel)
12
13     try:
14         request = backward_service_pb2.SleepCycleRequest()
15         response = stub.RunSleepCycle(request)
16
17         return {
18             "status": "success",
19             "num_schemas": response.num_schemas,
20             "duration_seconds": response.duration_seconds,
21         }
22     except grpc.RpcError as e:
23         # Retry on failure
24         raise self.retry(exc=e, countdown=60)
25
26 @celery_app.task
```

```
27 def backtest_strategy(start_date: str, end_date: str):
28     """Long-running backtesting job."""
29     # Call Backward service for historical simulation
30     # This can take hours, so runs in background
31     pass
```

6 Hybrid Training Pipeline

6.1 PyTorch Training (Phase 1)

```
1 # training/train_vivit.py
2 import torch
3 import torch.nn as nn
4 from torch.utils.data import DataLoader
5
6 class ViViTTrainer:
7     def __init__(self, model, device='cuda'):
8         self.model = model.to(device)
9         self.device = device
10        self.optimizer = torch.optim.AdamW(model.parameters(), lr=1
11            e-4)
12
13    def train_epoch(self, dataloader):
14        self.model.train()
15        total_loss = 0
16
17        for batch in dataloader:
18            videos, labels = batch
19            videos = videos.to(self.device)
20            labels = labels.to(self.device)
21
22            # Forward pass
23            outputs = self.model(videos)
24            loss = nn.CrossEntropyLoss()(outputs, labels)
25
26            # Backward pass
27            self.optimizer.zero_grad()
28            loss.backward()
29            self.optimizer.step()
```

```

30         total_loss += loss.item()
31
32     return total_loss / len(dataloader)
33
34     def export_to_onnx(self, path: str):
35         """Export trained model to ONNX for Rust inference."""
36         self.model.eval()
37         dummy_input = torch.randn(1, 16, 2, 60, 60).to(self.device)
38
39         torch.onnx.export(
40             self.model,
41             dummy_input,
42             path,
43             export_params=True,
44             opset_version=14,
45             do_constant_folding=True,
46             input_names=['video'],
47             output_names=['embedding'],
48             dynamic_axes={
49                 'video': {0: 'batch_size'},
50                 'embedding': {0: 'batch_size'}
51             }
52         )
53         print(f"Model exported to {path}")
54
55     # Usage
56     trainer = ViViTTrainer(model)
57     for epoch in range(100):
58         loss = trainer.train_epoch(train_loader)
59         print(f"Epoch {epoch}: Loss = {loss}")
60
61     # Export to ONNX
62     trainer.export_to_onnx("models/vivit.onnx")

```

6.2 Model Export & Validation

```

1 # training/validate_export.py
2 import torch
3 import onnx
4 import onnxruntime as ort
5 import numpy as np

```

```

6
7 def validate_onnx_export(pytorch_model, onnx_path):
8     """Validate that ONNX export matches PyTorch output."""
9
10    # Load ONNX model
11    onnx_model = onnx.load(onnx_path)
12    onnx.checker.check_model(onnx_model)
13
14    # Create ONNX Runtime session
15    ort_session = ort.InferenceSession(onnx_path)
16
17    # Test input
18    dummy_input = torch.randn(1, 16, 2, 60, 60)
19
20    # PyTorch inference
21    pytorch_model.eval()
22    with torch.no_grad():
23        pytorch_output = pytorch_model(dummy_input).numpy()
24
25    # ONNX inference
26    ort_inputs = {ort_session.get_inputs()[0].name: dummy_input.
27                  numpy()}
28    onnx_output = ort_session.run(None, ort_inputs)[0]
29
30    # Compare outputs
31    diff = np.abs(pytorch_output - onnx_output).max()
32    print(f"Max difference: {diff}")
33
34    if diff < 1e-5:
35        print("[OK] ONNX export validated successfully")
36    else:
37        print("[FAIL] ONNX export validation failed")
38
39    validate_onnx_export(model, "models/vivit.onnx")

```

7 Deployment Architecture

7.1 Docker Compose Setup

```

1 # docker-compose.yml

```

```
2 version: '3.8'
3
4 services:
5   gateway:
6     build: ./gateway
7     ports:
8       - "8000:8000"
9     environment:
10       - FORWARD_SERVICE_URL=forward:50051
11       - BACKWARD_SERVICE_URL=backward:50052
12     depends_on:
13       - forward
14       - backward
15       - redis
16
17   forward:
18     build: ./services/forward
19     ports:
20       - "50051:50051"
21     volumes:
22       - ./models:/models:ro
23     environment:
24       - RUST_LOG=info
25       - MODEL_PATH=/models/vivit.onnx
26     deploy:
27       resources:
28         limits:
29           cpus: '4'
30           memory: 2G
31
32   backward:
33     build: ./services/backward
34     ports:
35       - "50052:50052"
36     volumes:
37       - ./models:/models:ro
38       - ./data:/data
39     environment:
40       - RUST_LOG=info
41       - QDRANT_URL=http://qdrant:6333
42     depends_on:
```

```
43     - qdrant
44
45   qdrant:
46     image: qdrant/qdrant:latest
47     ports:
48       - "6333:6333"
49     volumes:
50       - qdrant_data:/qdrant/storage
51
52   redis:
53     image: redis:7-alpine
54     ports:
55       - "6379:6379"
56
57   celery_worker:
58     build: ./gateway
59     command: celery -A tasks worker --loglevel=info
60     depends_on:
61       - redis
62       - backward
63
64   volumes:
65     qdrant_data:
```

7.2 Kubernetes Deployment

```
1  # k8s/forward-deployment.yaml
2  apiVersion: apps/v1
3  kind: Deployment
4  metadata:
5    name: janus-forward
6  spec:
7    replicas: 3
8    selector:
9      matchLabels:
10        app: janus-forward
11    template:
12      metadata:
13        labels:
14          app: janus-forward
15      spec:
```

```
16     containers:
17     - name: forward
18       image: janus/forward:latest
19       ports:
20       - containerPort: 50051
21       resources:
22         requests:
23           memory: "1Gi"
24           cpu: "2"
25         limits:
26           memory: "2Gi"
27           cpu: "4"
28       env:
29       - name: RUST_LOG
30         value: "info"
31       - name: MODEL_PATH
32         value: "/models/vivit.onnx"
33       volumeMounts:
34       - name: models
35         mountPath: /models
36         readOnly: true
37       volumes:
38       - name: models
39         persistentVolumeClaim:
40           claimName: model-storage
41 ---
42 apiVersion: v1
43 kind: Service
44 metadata:
45   name: janus-forward
46 spec:
47   selector:
48     app: janus-forward
49   ports:
50   - protocol: TCP
51     port: 50051
52     targetPort: 50051
53   type: ClusterIP
```

8 Implementation Roadmap

8.1 Phase 1: Foundation (Weeks 1-4)

1. Week 1: Project Setup

- Initialize Rust workspace with cargo
- Set up CI/CD (GitHub Actions)
- Create protobuf definitions
- Generate gRPC stubs for Rust and Python

2. Week 2: Core Crates

- Implement `janus-core` with domain types
- Implement `janus-vision` GAF transformation
- Add comprehensive unit tests

3. Week 3: Model Integration

- Train ViViT model in PyTorch
- Export to ONNX and validate
- Integrate ONNX Runtime in Rust
- Benchmark inference latency

4. Week 4: Basic Services

- Implement Forward service skeleton
- Implement Backward service skeleton
- Set up FastAPI gateway
- Test end-to-end flow

8.2 Phase 2: Core Features (Weeks 5-8)

1. Week 5: LTN Implementation

- Implement Łukasiewicz t-norms
- Encode wash sale constraint
- Encode Almgren-Chriss constraint
- Encode VPIN constraint

2. Week 6: Decision Engine

- Implement dual pathway logic
- Add risk gating
- Add logic gating
- Implement cerebellar forward model

3. Week 7: Memory System

- Implement episodic buffer
- Implement prioritized replay
- Add schema consolidation

4. Week 8: Vector Database

- Set up Qdrant
- Implement schema storage
- Implement similarity search
- Add periodic pruning

8.3 Phase 3: Production Readiness (Weeks 9-12)

1. Week 9: Performance Optimization

- Profile and optimize hot paths
- Add SIMD optimizations
- Implement model quantization
- Benchmark end-to-end latency

2. Week 10: Safety & Testing

- Remove all panic!() calls
- Add comprehensive error handling
- Write integration tests
- Add property-based tests

3. Week 11: Monitoring & Observability

- Add Prometheus metrics
- Implement distributed tracing

- Set up alerting
- Create dashboards

4. Week 12: Deployment

- Create Docker images
- Write Kubernetes manifests
- Set up staging environment
- Run load tests

9 Complete Implementation Checklist

sec:checklist

9.1 Infrastructure Setup

- Create Rust workspace (`Cargo.toml`)
- Set up monorepo structure (services + crates)
- Configure CI/CD pipeline
- Set up Docker build system
- Configure linting (clippy, rustfmt)

9.2 Shared Crates

- **janus-core**: Domain types, errors, utilities
- **janus-vision**: GAF, ViViT, image processing
- **janus-logic**: LTN, constraint evaluation
- **janus-memory**: Replay buffer, schemas
- **janus-execution**: Order execution, market simulation
- **janus-proto**: Protobuf definitions and gRPC stubs

9.3 Forward Service (Rust)

- gRPC server setup with Tonic
- GAF transformation pipeline
- ONNX model loading and inference
- LTN constraint evaluation
- Multimodal fusion
- Decision engine with dual pathways
- Metrics export (Prometheus)
- Health check endpoint
- Graceful shutdown

9.4 Backward Service (Rust)

- gRPC server setup
- Episodic buffer implementation
- Prioritized replay sampling
- TD-error computation
- Schema consolidation
- UMAP integration
- Qdrant client and sync
- Sleep cycle orchestration
- Parallel batch processing

9.5 Gateway Service (Python)

- FastAPI application setup
- gRPC client stubs
- REST endpoints (/decision, /sleep-cycle, etc.)
- JWT authentication

- Rate limiting middleware
- Celery worker setup
- Background job management
- Health check aggregation
- Metrics aggregation

9.6 Training Pipeline (Python)

- PyTorch model definitions
- Data loaders and preprocessing
- Training loops
- ONNX export scripts
- Export validation
- Model versioning

9.7 Testing

- Unit tests for all crates
- Integration tests for services
- End-to-end tests
- Load tests
- Chaos testing

9.8 Deployment

- Docker images for all services
- Docker Compose for local dev
- Kubernetes manifests
- Helm charts
- Staging environment
- Production environment

10 Conclusion

This Rust-first implementation strategy for JANUS provides:

1. **Maximum Performance:** Rust's zero-cost abstractions and no GIL enable <100ms latency for Forward service
2. **Type Safety:** Compile-time guarantees prevent entire classes of bugs
3. **Memory Safety:** No null pointer dereferences, no data races
4. **Developer Experience:** Python gateway maintains ease of use for API consumers
5. **Ecosystem Compatibility:** Hybrid approach leverages best tools from both ecosystems

Next Steps:

- Begin with Phase 1 foundation (weeks 1-4)
- Train initial models in PyTorch and export to ONNX
- Implement Forward service with ONNX Runtime
- Gradually migrate components to pure Rust (Candle/Burn)
- Monitor performance and iterate

The architecture is designed for **incremental migration**, allowing you to start with ONNX inference in Rust while keeping training in PyTorch, then gradually move to a fully Rust-native stack as the ecosystem matures.

References

- [1] Jordan Smith, "JANUS Forward: Wake State Logic Trading Algorithm," 2025.
- [2] Jordan Smith, "JANUS Backward: Sleep State Memory Management," 2025.
- [3] Laurent Mazare, "tch-rs: Rust bindings for PyTorch," GitHub, 2024.
- [4] Microsoft, "ONNX Runtime: Cross-platform ML inference," 2024.
- [5] Hugging Face, "Candle: Minimalist ML framework in Rust," 2024.
- [6] Burn Contributors, "Burn: Deep learning framework in Rust," 2024.
- [7] Tokio Contributors, "Tonic: A native gRPC client & server implementation," 2024.
- [8] Sebastián Ramírez, "FastAPI: Modern, fast web framework for Python," 2024.
- [9] Rayon Contributors, "Rayon: Data parallelism library for Rust," 2024.
- [10] ndarray Contributors, "ndarray: N-dimensional arrays for Rust," 2024.
- [11] Qdrant Team, "Qdrant: Vector database for ML applications," 2024.