

# SPL Integration Guide v3.1

**\*\*Step-by-Step Implementation for Existing LLM Agents\*\***

**\*\*Author:\*\*** Pamela Cuce  
**\*\*Email:\*\*** pamela@dasein.works  
**\*\*Phone:\*\*** 914-240-3564  
**\*\*Version:\*\*** 3.1 (Corrected)  
**\*\*Date:\*\*** December 3, 2025

---

**## Integration Overview**

SPL integrates as a **\*\*wrapper layer\*\*** around existing LLM calls. No API changes needed.

**BEFORE (Direct LLM):**

request → LLM → response

**AFTER (SPL Wrapper):**

request → Layer 0 → Layer 1 → Layer 2 (LLM) → response

**\*\*Integration time:\*\*** 30 minutes

**\*\*LOC change:\*\*** 50-100 lines

**\*\*Breaking changes:\*\*** None

---

**## Step 1: Install Dependencies**

Python 3.8+

Install Anthropic SDK for LLM access

pip install anthropic

Optional: for multi-agent state sharing

pip install redis

--

## Step 2: Implement Layer 0 (Validation)

### Purpose

Layer 0 performs structural validation - rejecting invalid requests before they reach expensive layers.

### Code with Full Documentation

### **Layer 0: Structural Validation (Free, Fast)**

**Purpose: Reject invalid inputs immediately, no LLM calls**

def layer\_0\_validate(request):

"""

Layer 0: Structural validation

Checks if request has required format/fields

Returns: {'halt': True/False, 'error': str or None}

Validations performed:

1. Check request has 'user\_id' (for authentication)
2. Check request has 'content' (for processing)

### 3. Check content is reasonable length

Cost: \$0 (no API calls)

Speed: <1ms (just dict/type checks)

.....

# Check 1: Does request have user\_id?

# This is required for:

# - Permission checking (is user allowed?)

# - Rate limiting (has user exceeded quota?)

# - Audit logging (who made this request?)

```
if not request.get('user_id'):
```

```
    return {'halt': True, 'error': 'Missing user_id'}
```

# Check 2: Does request have content?

# This is what we're going to categorize/analyze

```
if not request.get('content'):
```

```
    return {'halt': True, 'error': 'Missing content'}
```

# Check 3: Is content reasonable length?

# Too short (< 10 chars) = probably spam/invalid

# Too long (> 100k chars) = probably spam/DoS attack

```
if len(request['content']) > 100000:
```

```
return {'halt': True, 'error': 'Content too long'}
```

  

```
if len(request['content']) < 10:
```

```
    return {'halt': True, 'error': 'Content too short'}
```

  

```
# All checks passed - safe to continue
```

```
return {'halt': False}
```

---

```
## Step 3: Implement Layer 1 (Patterns)
```

```
### Purpose
```

Layer 1 uses pattern matching and rules to suppress expensive LLM calls.

```
### Code with Full Documentation
```

## **Layer 1: Pattern Matching & Rules (Cheap, Medium Speed)**

**Purpose: Match against known patterns before calling LLM**

```
import re
```

### Pre-defined high-accuracy patterns

```
LAYER_1_PATTERNS = {  
    'urgent': {  
        'regex': r'urgent|asap|emergency|immediately', # Pattern to match  
        'category': 'urgent', # Category if matched  
        'confidence': 0.95 # 95% accuracy  
    },  
    'billing': {  
        'regex': r'invoice|payment|bill|receipt|charge',  
        'category': 'billing',  
        'confidence': 0.92  
    },  
    'spam': {  
        'regex': r'unsubscribe|viagra|lottery|winner|claim',  
        'category': 'spam',  
    }  
}  
  
def layer_1_match(content, user_patterns=None):  
    """
```

### Layer 1: Pattern matching

Tries to match content against known patterns

Returns: {'matched': True/False, 'category': str, 'confidence': float}

### Strategy:

1. Combine pre-defined patterns with user-provided patterns
2. For each pattern, try regex match against content
3. If match found AND confidence >= 0.85: suppress Layer 2
4. If no match or low confidence: escalate to Layer 2

Cost: \$0 (just regex matching)

Speed: <5ms (even with 100+ patterns)

Accuracy: varies by pattern (0.85-0.98)

.....

# Merge pre-defined and user patterns

```
patterns = {**LAYER_1_PATTERNS, **(user_patterns or {})}
```

```
content_lower = content.lower()
```

# Try each pattern

```
for pattern_name, pattern_config in patterns.items():
```

# regex: pattern to search for

# category: category if pattern matches

# confidence: how confident are we in this pattern?

```
regex = pattern_config['regex']
```

# Try to match pattern

# re.search() finds pattern anywhere in string

```
if re.search(regex, content_lower):
```

```
    confidence = pattern_config['confidence']
```

```
# Is confidence high enough to skip expensive LLM?  
# Threshold: 0.85 (tunable - change for your use case)  
  
if confidence >= 0.85:  
  
    # Yes! Suppress Layer 2  
  
    return {  
  
        'matched': True,  
  
        'category': pattern_config['category'],  
  
        'confidence': confidence,  
  
        'cost': 0.0  
  
    }  
  
  
# No patterns matched (or all were below threshold)  
return {'matched': False}
```

---

```
## Step 4: Implement Layer 2 (LLM)
```

### Purpose

Layer 2 calls LLM for complex reasoning when Layers 0-1 can't decide.

### Code with Full Documentation

## **Layer 2: LLM Reasoning (Expensive, Slow, Powerful)**

**Purpose: Call LLM for complex reasoning when simpler layers fail**

import anthropic

```
def layer_2_llm(content, model='claude-3-5-sonnet-20241022'):
```

```
"""
```

Layer 2: LLM reasoning

Calls LLM when Layers 0-1 couldn't decide

Returns: {'category': str, 'confidence': float, 'cost': float}

When called:

- Layer 0 passed validation
- Layer 1 found no high-confidence pattern match
- Must use full LLM reasoning

Cost: ~\$0.01 per call (expensive!)

Speed: 100-300ms (medium)

Accuracy: 90-95% typical

Setup:

1. Create Anthropic client
2. Build prompt for categorization
3. Call LLM API

4. Extract category from response

5. Return with cost metadata

.....

```
# Initialize Anthropic client
```

```
# Uses ANTHROPIC_API_KEY environment variable by default
```

```
client = anthropic.Anthropic()
```

```
# Build prompt for LLM
```

```
# This tells LLM what to do
```

```
prompt = f"""Categorize this email into one category: urgent, billing, spam, or other.
```

```
Email: {content}
```

```
Respond with ONLY the category name."""
```

```
# Call LLM API
```

```
# model: which LLM model to use
```

```
# max_tokens: limit response length
```

```
# temperature: 0=deterministic, 1=creative (use 0 for consistency)
```

```
message = client.messages.create(
```

```
    model=model,
```

```
    max_tokens=50, # We only need 1 word response
```

```
    temperature=0, # Deterministic (no randomness)
```

```
    messages=[
```

```
{  
    "role": "user",  
    "content": prompt  
}  
]  
)  
  
# Extract category from response  
  
# LLM returns structured message, get text content  
category = message.content.text.strip().lower()  
  
  
# Validate category (in case LLM returns unexpected response)  
valid_categories = ['urgent', 'billing', 'spam', 'other']  
if category not in valid_categories:  
    # Unknown category - return as 'other'  
    category = 'other'  
  
  
# Return with cost information  
# Typical cost: $0.01 per request  
return {  
    'category': category,  
    'confidence': 0.92, # LLM usually ~92% confident  
    'cost': 0.01,      # API cost per request
```

```
'model': model      # Track which model was used  
}  
  
---
```

## ## Step 5: Wire Layers Together

### ### Purpose

Orchestrate the three layers into a complete pipeline.

### ### Code with Full Documentation

#### **Main SPL Pipeline: Orchestrate all three layers**

```
def spl_process(request):
```

```
    """
```

Main SPL pipeline: execute all three layers

Args:

    request: dict with 'user\_id' and 'content'

Returns:

    dict with result and metadata

Execution flow:

## 1. Layer 0: Validate structure

- If invalid → REJECT (return immediately)
- If valid → continue to Layer 1

## 2. Layer 1: Pattern matching

- If high-confidence match → RETURN RESULT (skip expensive Layer 2!)
- If no match → continue to Layer 2

## 3. Layer 2: LLM reasoning

- Call LLM (expensive!)
- Return result

Key insight:

- Most requests (80-95%) stop at Layer 0-1
- Only 5-20% reach expensive Layer 2
- This is how we achieve 10-50x cost reduction

.....

```
# ===== LAYER 0: VALIDATION =====  
  
# Check if request is valid  
validation = layer_0_validate(request)  
if validation['halt']:  
    # Invalid request - reject immediately  
    # Cost: $0  
    return {  
        'status': 'rejected',  
        'reason': validation['error'],  
        'cost': 0.0,  
        'layer': 0
```

```
}

# ===== LAYER 1: PATTERN MATCHING =====

# Try to match against known patterns
pattern_result = layer_1_match(request['content'])

if pattern_result['matched']:
    # Pattern matched - suppress Layer 2 (skip expensive LLM!)
    # Cost: $0 (just pattern matching)
    return {
        'status': 'success',
        'category': pattern_result['category'],
        'confidence': pattern_result['confidence'],
        'cost': pattern_result['cost'],
        'layer': 1,
        'method': 'pattern'
    }

# ===== LAYER 2: LLM =====

# No pattern matched - must use LLM
llm_result = layer_2_llm(request['content'])

# Cost: $0.01 (expensive!)
return {
    'status': 'success',
    'category': llm_result['category'],
    'confidence': llm_result['confidence'],
    'cost': llm_result['cost'],
    'layer': 2,
    'method': 'llm'
```

```
}
```

```
--
```

## ## Step 6: Add Cost Tracking

### ### Purpose

Monitor and report costs for optimization.

### ### Code with Full Documentation

#### **Cost Tracking: Monitor spending and efficiency**

class CostTracker:

```
"""
```

Track costs and efficiency metrics

Attributes:

total\_cost: cumulative cost in dollars

llm\_calls: count of Layer 2 LLM calls

pattern\_hits: count of Layer 1 pattern matches

rejections: count of Layer 0 rejections

Metrics calculated:

- cost per request
- suppression rate (% avoiding LLM)
- cost reduction factor

```
"""
```

```
def __init__(self):
    """Initialize cost tracking"""
    self.total_cost = 0.0
    self.llm_calls = 0      # Expensive calls
    self.pattern_hits = 0   # Cheap hits
    self.rejections = 0     # Free rejections
```

```
def record(self, result):
```

```
    """
```

```
    Record one request result
```

Args:

result: dict from spl\_process() with layer and cost

Updates:

- total\_cost: add this request's cost
- layer counters: track which layer made decision

```
    """
```

```
# Add to total cost
```

```
    self.total_cost += result['cost']
```

```
# Track by layer
```

```
if result['status'] == 'rejected':
```

```
    self.rejections += 1
```

```
elif result['layer'] == 1:
```

```
    self.pattern_hits += 1
```

```
elif result['layer'] == 2:
```

```
    self.llm_calls += 1
```

```
def report(self):
    """
    Generate cost report with key metrics

```

Returns:

dict with:

- total\_cost: dollars spent
- cost\_per\_request: average cost per request
- suppression\_rate: % of requests avoiding LLM
- cost\_reduction\_factor: X times cheaper than baseline

Example output:

```
{
    'total_cost': $5.00,
    'cost_per_request': $0.001,
    'suppression_rate': 95%,
    'cost_reduction_factor': 10x
}
```

```
total_requests = self.llm_calls + self.pattern_hits + self.rejections
```

```
# Calculate metrics
cost_per_request = self.total_cost / total_requests if total_requests > 0 else 0
suppression_rate = (self.pattern_hits + self.rejections) / total_requests if
total_requests > 0 else 0
```

```
# Baseline: if all requests used LLM
baseline_cost = total_requests * 0.01
```

```

cost_reduction = baseline_cost / self.total_cost if self.total_cost > 0 else 0

return {
    'total_cost': self.total_cost,
    'total_requests': total_requests,
    'llm_calls': self.llm_calls,
    'pattern_hits': self.pattern_hits,
    'rejections': self.rejections,
    'cost_per_request': cost_per_request,
    'suppression_rate': suppression_rate,
    'cost_reduction_factor': cost_reduction
}

```

### **Usage example:**

```

tracker = CostTracker()

for request in requests:
    result = spl_process(request)
    tracker.record(result)

print(tracker.report())

```

**Output: {'total\_cost': 2.50, 'suppression\_rate': 0.75, 'cost\_reduction\_factor': 4.0x}**

---

```

## Step 7: Integrate with Existing Agent

### Drop-in Replacement Pattern

**Before (Direct LLM):**

def categorize_email(email):

```

```
"""Original agent - calls LLM directly"""

# Direct LLM call - costs $0.01 per email
```

```
response = llm.complete(
    prompt=f"Categorize: {email['subject']}",
    max_tokens=100
)

return response['category']
```

```
**After (SPL Wrapper):**
```

```
def categorize_email_with_spl(email):
    """Agent with SPL wrapper - costs $0.001 on average"""

    # Prepare request for SPL
    request = {
        'user_id': email['from'],
        'content': email['subject']
    }

    # Process through SPL pipeline
    # Result: 80-95% bypass expensive LLM!
    result = spl_process(request)

    if result['status'] == 'rejected':
        return 'error'

    return result['category']
```

## That's it! Drop-in replacement with 10x cost savings

--

```
## Step 8: Test & Validate
```

```
### Comprehensive Testing
```

```
def test_spl_integration():
```

```
    """
```

```
Test SPL with sample emails
```

```
Verifies:
```

- Correct categorization
- Cost optimization
- Pattern learning

```
    """
```

```
# Test cases: (input, expected_category)
```

```
test_cases = [
```

```
    ('URGENT: Meeting at 3pm', 'urgent'),  
    ('Invoice #12345', 'billing'),  
    ('You won $1000!', 'spam'),  
    ('Team lunch tomorrow', 'other'),
```

```
]
```

```
cost_tracker = CostTracker()
```

```
for content, expected_category in test_cases:
```

```
    request = {'user_id': 'test', 'content': content}  
    result = spl_process(request)
```

```
# Verify categorization
```

```
    print(f"Input: {content}")  
    print(f" Expected: {expected_category}")  
    print(f" Got: {result['category']}")  
    print(f" Layer: {result['layer']}")  
    print(f" Cost: ${result['cost']:.4f}")  
    print()
```

```
cost_tracker.record(result)
```

```
# Print summary
print("Overall Report:")
print(cost_tracker.report())
```

---

```
## Step 9: Deploy to Production
```

```
### Configuration Management
```

## **config.py - Production configuration**

```
SPL_CONFIG = {
    # Layer 0: Validation
    'layer_0': {
        'enabled': True,
        'rate_limit': 100 # requests/minute per user
    },
    # Layer 1: Pattern matching
    'layer_1': {
        'enabled': True,
        'confidence_threshold': 0.85,
        'patterns_file': 'patterns.json' # Load patterns from file
    },
    # Layer 2: LLM
    'layer_2': {
        'enabled': True,
        'model': 'claude-3-5-sonnet-20241022',
        'timeout_ms': 5000, # Timeout protection
        'max_retries': 3 # Retry on API failure
    },
    # Monitoring
    'monitoring': {
        'enabled': True,
        'log_level': 'INFO',
        'metrics_enabled': True
    }
}
```

```
}
```

```
}
```

### ### Error Handling with Fallback

```
def spl_process_safe(request):
```

```
"""
```

```
SPL with comprehensive error handling
```

Handles:

- API failures (fallback to Layer 1)
- Timeouts (use cached result)
- Rate limits (graceful degradation)
- Malformed responses (use 'other')

```
"""
```

```
import logging
```

```
logger = logging.getLogger(__name__)
```

```
try:
```

```
    # Layer 0: Validation
```

```
    validation = layer_0_validate(request)
```

```
    if validation['halt']:
```

```
        return {'status': 'rejected', 'reason': validation['error'], 'cost': 0.0}
```

```
    # Layer 1: Pattern matching
```

```
    pattern_result = layer_1_match(request['content'])
```

```
    if pattern_result['matched']:
```

```
        return {'status': 'success', 'category': pattern_result['category'], 'cost': 0.0, 'layer': 1}
```

```
    # Layer 2: LLM (with error handling)
```

```
    try:
```

```
        llm_result = layer_2_llm(request['content'])
```

```
        return {'status': 'success', 'category': llm_result['category'], 'cost': 0.01, 'layer': 2}
```

```
    except anthropic.APIError as e:
```

```

# API failed - fallback to Layer 1 result
logger.error(f"LLM error: {e}, falling back to Layer 1")
return {'status': 'fallback', 'error': str(e), 'cost': 0.0}

except anthropic.Timeout as e:
    # LLM timed out - use pattern fallback
    logger.warning(f"LLM timeout: {e}, using patterns")
    return {'status': 'timeout', 'error': str(e), 'cost': 0.0}

except Exception as e:
    # Unexpected error - log and return error
    logger.error(f"Unexpected error: {e}")
    return {'status': 'error', 'error': str(e), 'cost': 0.0}

```

---

**## Step 10: Monitor & Optimize**

### Metrics Dashboard

## Logging for monitoring

```

import logging
from datetime import datetime

logging.basicConfig(level=logging.INFO)
logger = logging.getLogger(name)

def log_result(result):
    """Log execution for monitoring"""

    # Log key metrics for later analysis
    logger.info(f"SPL Result: layer={result['layer']}, "
               f"cost=${{result['cost']:.4f}}, "
               f"status={{result['status']}}, "
               f"category={{result.get('category', 'N/A')}}")

```

### In your main loop:

```

for request in incoming_requests:
    result = spl_process_safe(request)
    log_result(result)

```

### ### Cost Alerts

```
def check_budget(tracker, daily_budget=100.0):
```

....

Alert if daily budget exceeded

Args:

tracker: CostTracker instance

daily\_budget: max allowed spend per day

Triggers cost-saving measures if over budget:

- Increase confidence threshold (fewer LLM calls)
- Reduce timeout (fallback faster)
- Add emergency patterns

....

```
if tracker.total_cost > daily_budget:
```

# Budget exceeded!

```
logger.warning(f"Budget alert: ${tracker.total_cost:.2f} "
               f"of ${daily_budget:.2f} spent")
```

# Trigger cost-saving measures

# (implementation depends on your needs)

--

### ## Performance Targets

**\*\*Week 1 (Baseline):\*\***

- Pattern hit rate: 50%
- Cost reduction: 2x
- Action: Add more patterns based on data

**\*\*Week 2 (Tuning):\*\***

- Pattern hit rate: 70%
- Cost reduction: 5x
- Action: Lower confidence threshold, add learned patterns

**\*\*Week 3+ (Optimization):\*\***

- Pattern hit rate: 85%+
- Cost reduction: 10-15x
- Action: Continuous pattern optimization

---

## ## Troubleshooting

### ### Problem: Layer 1 patterns not matching

**\*\*Solution:\*\***

- Review pattern regex for bugs
- Lower confidence threshold temporarily
- Add more comprehensive patterns
- Check content preprocessing (case sensitivity, whitespace)

### ### Problem: Layer 2 calls still too expensive

**\*\*Solution:\*\***

- Increase pattern coverage (more patterns = fewer LLM calls)
- Implement result caching (avoid reprocessing identical content)
- Use cheaper LLM model (trade accuracy for cost)
- Batch requests (process in bulk for better pricing)

### Problem: Rate limiting too aggressive

\*\*Solution:\*\*

- Increase rate\_limit threshold
- Implement user-specific limits
- Add priority queuing for important users

### Problem: Timeouts on Layer 2

\*\*Solution:\*\*

- Add timeout protection
- Increase timeout\_ms (allow more time)
- Fallback to Layer 1 results on timeout

---

## ## Next Steps

1.  Implement Steps 1-10
2.  Deploy to staging environment
3.  Run cost analysis for 1 week
4.  Document results and learnings
5.  Deploy to production
6.  Monitor cost reduction in real-time
7.  Optimize patterns based on live data
8.  Scale to multi-agent network (see SPL-MultiAgent-v3.1.md)

--  
\*Pamela Cuce | pamcuce@alum.vassar.edu | Version 3.1 (Corrected) | December 3, 2025\*