

Case Studies: Rule-Based NLP with spaCy

Case Study 1: Medical Document Processing System for Clinical Reports

Background

Organization: HealthTech Solutions, a healthcare analytics company

Challenge: Process thousands of unstructured clinical reports daily to extract key medical information including diseases, medications, symptoms, and treatment procedures. Pre-trained NLP models miss hospital-specific terminology and medication names.

Problem Statement

The company receives 5,000+ clinical reports daily containing:

- Custom medication codes (e.g., "MED-2045", "RX-COVID-VAC-01")
- Hospital-specific disease terminology (e.g., "T2DM" for Type 2 Diabetes)
- Local treatment protocols and procedure names
- COVID-19 related terms and variants
- Adverse drug reaction descriptions

Challenges:

1. Generic NLP models don't recognize hospital-specific codes
2. New medications and treatments emerge frequently
3. Need 100% accuracy for regulatory compliance
4. Must process reports in real-time (< 2 seconds per document)
5. Limited labeled training data available

Solution: Rule-Based NLP with spaCy

Implementation Components

1. PhraseMatcher for Medical Terms

```
python
```

```

import spacy
from spacy.matcher import PhraseMatcher

# Load spaCy model
nlp = spacy.load("en_core_web_sm")

# Define medical term dictionaries
disease_terms = [
    "Type 2 Diabetes", "T2DM", "Diabetes Mellitus Type 2",
    "Hypertension", "High Blood Pressure", "HTN",
    "COVID-19", "SARS-CoV-2", "Coronavirus Disease",
    "Myocardial Infarction", "Heart Attack", "MI",
    "Chronic Kidney Disease", "CKD", "Renal Failure"
]

medication_terms = [
    "Metformin", "Insulin Glargine", "Lantus",
    "Lisinopril", "ACE Inhibitor",
    "Remdesivir", "Paxlovid", "COVID-19 Vaccine",
    "Aspirin", "Acetylsalicylic Acid"
]

# Create PhraseMatcher
matcher = PhraseMatcher(nlp.vocab, attr="LOWER")
patterns_disease = [nlp.make_doc(text) for text in disease_terms]
patterns_meds = [nlp.make_doc(text) for text in medication_terms]

matcher.add("DISEASE", patterns_disease)
matcher.add("MEDICATION", patterns_meds)

# Process document
doc = nlp("Patient diagnosed with T2DM, prescribed Metformin 500mg twice daily.")
matches = matcher(doc)

for match_id, start, end in matches:
    span = doc[start:end]
    label = nlp.vocab.strings[match_id]
    print(f"\n{label}: {span.text}\n")

```

Output:

DISEASE: T2DM

MEDICATION: Metformin

2. EntityRuler for Custom Medical Codes

python


```
print("Extracted Medical Entities:")
print("-" * 60)
for ent in doc.ents:
    print(f'{ent.label_}: {ent.text}')
```

Output:

Extracted Medical Entities:

```
COVID_TERM      | COVID-19 positive
DISEASE        | T2DM
DISEASE_CODE    | ICD-10:E11.9
DISEASE         | HTN
MEDICATION_CODE | MED-2045
MEDICATION      | Metformin
DOSAGE          | 500mg
MEDICATION_CODE | RX-COVID-VAC-01
ADVERSEREACTION | side effect
MEDICATION      | Metformin
DOSAGE          | 500mg
```

3. Token Matcher for Symptom Detection

python

```

from spacy.matcher import Matcher

matcher = Matcher(nlp.vocab)

# Pattern for symptoms
symptom_patterns = [
    [{"LOWER": "fever"}, {"LOWER": {"IN": ["of", ">"]}}, {"LIKE_NUM": True}],
    [{"LOWER": "difficulty"}, {"LOWER": "breathing"}],
    [{"LOWER": "shortness"}, {"LOWER": "of"}, {"LOWER": "breath"}],
    [{"LOWER": "chest"}, {"LOWER": "pain"}],
    [{"LOWER": "persistent"}, {"LOWER": "cough"}],
]

for pattern in symptom_patterns:
    matcher.add("SYMPTOM", [pattern])

# Process
text = "Patient presents with fever of 102.5°F, difficulty breathing, and persistent cough."
doc = nlp(text)
matches = matcher(doc)

print("\nSymptoms Detected:")
for match_id, start, end in matches:
    print(f"- {doc[start:end].text}")

```

Output:

Symptoms Detected:

- fever of 102.5
- difficulty breathing
- persistent cough

Results & Impact

Quantitative Outcomes:

- **Accuracy:** 99.2% for medication extraction (vs 78% with generic NLP)
- **Processing Speed:** 1.3 seconds per document (met < 2 second requirement)
- **Coverage:** Successfully extracted 15,000+ unique medical entities daily
- **Cost Savings:** \$2M annually by automating manual review
- **Compliance:** 100% regulatory compliance achieved

Qualitative Benefits:

- Real-time alerts for adverse drug reactions
- Automated patient risk scoring
- Streamlined insurance claim processing
- Improved clinical decision support
- Faster response to emerging health threats (COVID-19 variants)

Key Learnings

1. PhraseMatcher Advantages:

- Perfect for large medical terminology dictionaries (10,000+ terms)
- Case-insensitive matching handles variations (T2DM, t2dm)
- Fast performance even with extensive term lists

2. EntityRuler Flexibility:

- Regex patterns caught hospital-specific codes
- Easy updates when new medications approved
- No retraining required for new patterns

3. Hybrid Approach:

- Combined rule-based extraction with spaCy's pre-trained NER
- Rules handled domain-specific terms
- ML model captured general medical concepts

4. Maintenance:

- Monthly updates to phrase lists (new medications)
- Version control for pattern changes
- A/B testing for pattern effectiveness

Future Enhancements

- Integration with UMLS (Unified Medical Language System) terminology
 - Multi-language support for international reports
 - Active learning to suggest new patterns from unmatched entities
 - API for real-time entity extraction
-

Case Study 2: Financial Compliance Monitoring System

Background

Organization: Global Finance Corp, a multinational investment bank

Challenge: Monitor internal communications (emails, chat, reports) for regulatory compliance violations, financial fraud indicators, and insider trading patterns across 50,000+ daily messages.

Problem Statement

The compliance team must detect:

- Insider trading language ("buy before announcement", "confidential merger")
- Market manipulation terms ("pump and dump", "coordinated buying")
- Regulatory violations (GDPR breaches, data sharing violations)
- Company-specific policy violations (trading windows, disclosure requirements)
- Money laundering indicators (unusual transaction patterns)

Challenges:

1. Financial jargon and abbreviations (FICC, OTC, AML, KYC)
2. Context-dependent violations (legitimate vs suspicious communications)
3. New schemes and terminology emerge constantly
4. False positives must be minimized (alert fatigue)
5. Multiple regulatory frameworks (SEC, FCA, GDPR)
6. Real-time monitoring required

Solution: Rule-Based NLP with spaCy

Implementation Components

1. PhraseMatcher for Compliance Terms

```
python
```

```
import spacy
from spacy.matcher import PhraseMatcher

nlp = spacy.load("en_core_web_lg")

# Define compliance violation categories
insider_trading_terms = [
    "buy before announcement",
    "sell before news",
    "confidential merger information",
    "unreleased earnings",
    "material nonpublic information",
    "MNPI",
    "tip from insider",
    "quiet period violation",
    "trading window closed",
    "blackout period"
]

market_manipulation_terms = [
    "pump and dump",
    "coordinated buying",
    "artificial price inflation",
    "wash trading",
    "spoofing",
    "layering orders",
    "painting the tape",
    "front running"
]

money_laundering_terms = [
    "structuring deposits",
    "smurfing",
    "shell company transfer",
    "offshore account",
    "cash intensive business",
    "high risk jurisdiction",
    "beneficial owner concealment",
    "unusual wire transfer pattern"
]

gdprViolation_terms = [
    "personal data sharing",
```

```

"customer information disclosed",
"privacy breach",
"consent not obtained",
"data protection violation",
"right to be forgotten ignored",
"cross-border data transfer"
]

# Create PhraseMatcher
compliance_matcher = PhraseMatcher(nlp.vocab, attr="LOWER")

# Add patterns
compliance_matcher.add("INSIDER_TRADING", [nlp.make_doc(term) for term in insider_trading_terms])
compliance_matcher.add("MARKET_MANIPULATION", [nlp.make_doc(term) for term in market_manipulation_terms])
compliance_matcher.add("MONEY_LAUNDERING", [nlp.make_doc(term) for term in money_laundering_terms])
compliance_matcher.add("GDPR_VIOLATION", [nlp.make_doc(term) for term in gdprViolation_terms])

# Sample email monitoring
email_text = """
From: trader@globalfinance.com
To: colleague@globalfinance.com
Subject: Quick heads up

Hey - just got confidential merger information from our client.
The deal closes next week. You might want to buy before announcement.
Keep this quiet, we're in the blackout period.

"""

doc = nlp(email_text)
matches = compliance_matcher(doc)

print("⚠️ COMPLIANCE ALERTS DETECTED:")
print("=" * 70)
for match_id, start, end in matches:
    violation_type = nlp.vocab.strings[match_id]
    matched_text = doc[start:end].text
    print(f"\n⚠️ Violation Type: {violation_type}")
    print(f" Matched Text: '{matched_text}'")
    print(f" Risk Level: HIGH")

```

Output:

 COMPLIANCE ALERTS DETECTED:

⚠️ Violation Type: INSIDER_TRADING

Matched Text: 'confidential merger information'

Risk Level: HIGH

⚠️ Violation Type: INSIDER_TRADING

Matched Text: 'buy before announcement'

Risk Level: HIGH

⚠️ Violation Type: INSIDER_TRADING

Matched Text: 'blackout period'

Risk Level: HIGH

2. EntityRuler for Financial Instruments & Entities

```
python
```

```

from spacy.pipeline import EntityRuler

ruler = nlp.add_pipe("entity_ruler", before="ner")

patterns = [
    # Stock ticker patterns
    {"label": "STOCK_TICKER", "pattern": [{"TEXT": {"REGEX": "^[A-Z]{1,5}$"}, {"LOWER": "stock"}}],},
    {"label": "STOCK_TICKER", "pattern": [{"IS_UPPER": True, "LENGTH": {">=": 1, "<=": 5}}, {"LOWER": {"IN": ["sh"}]}]}

    # ISIN codes (International Securities Identification Number)
    {"label": "ISIN_CODE", "pattern": [{"TEXT": {"REGEX": "^[A-Z]{2}[A-Z0-9]{9}[0-9]$"}},]}

    # Account numbers
    {"label": "ACCOUNT_NUMBER", "pattern": [{"TEXT": {"REGEX": "^\u00c1CC-\d{8}$"}},]}, 
    {"label": "ACCOUNT_NUMBER", "pattern": [{"LOWER": "account"}, {"TEXT": {"REGEX": "^\u00d1d{10,12}$"}},]}

    # Transaction IDs
    {"label": "TRANSACTION_ID", "pattern": [{"TEXT": {"REGEX": "^\u00c1TXN-[A-Z0-9]{12}$"}},]}

    # Currency amounts (large suspicious amounts)
    {"label": "LARGE_AMOUNT", "pattern": [{"TEXT": "$"}, {"TEXT": {"REGEX": "^\u00d1[1-9]\d{6,}$"}},], "# IM+"}

    # Compliance terms
    {"label": "COMPLIANCE_TERM", "pattern": "material nonpublic information"}, 
    {"label": "COMPLIANCE_TERM", "pattern": "Chinese Wall"}, 
    {"label": "COMPLIANCE_TERM", "pattern": "Information Barrier"}, 
    {"label": "COMPLIANCE_TERM", "pattern": [{"UPPER": "MNPI"}]}, 
    {"label": "COMPLIANCE_TERM", "pattern": [{"UPPER": "AML"}], "# Anti-Money Laundering"}, 
    {"label": "COMPLIANCE_TERM", "pattern": [{"UPPER": "KYC"}], "# Know Your Customer"})

    # Suspicious action patterns
    {"label": "SUSPICIOUS_ACTION", "pattern": [{"LOWER": "urgent"}, {"LOWER": "wire"}, {"LOWER": "transfer"}]}, 
    {"label": "SUSPICIOUS_ACTION", "pattern": [{"LOWER": "delete"}, {"LOWER": "all"}, {"LOWER": "emails"}]}, 
    {"label": "SUSPICIOUS_ACTION", "pattern": [{"LOWER": "destroy"}, {"LOWER": "documents"}]}, 
    {"label": "SUSPICIOUS_ACTION", "pattern": [{"LOWER": "cash"}, {"LOWER": "transaction"}]}, 

]

ruler.add_patterns(patterns)

    # Monitor suspicious communication
    communication = """
    URGENT: Need to execute wire transfer of $15000000 to account ACC-98765432
    in Cayman Islands today. ISIN US0378331005 position must be liquidated.

```

Use TXN-ABC123XYZ789. Delete all emails after completion.

This involves MNPI from the board meeting. Chinese Wall protocols ignored.

.....

```
doc = nlp(communication)
```

```
print("\n FINANCIAL ENTITIES EXTRACTED:")
```

```
print("=" * 70)
```

```
for ent in doc.ents:
```

```
    print(f"<{ent.label_}:25} | {ent.text}")
```

```
# Calculate risk score
```

```
risk_score = 0
```

```
risk_factors = []
```

```
for ent in doc.ents:
```

```
    if ent.label_ == "SUSPICIOUS_ACTION":
```

```
        risk_score += 30
```

```
        risk_factors.append(f"Suspicious action: {ent.text}")
```

```
    elif ent.label_ == "COMPLIANCE_TERM":
```

```
        risk_score += 25
```

```
        risk_factors.append(f"Compliance term: {ent.text}")
```

```
    elif ent.label_ == "LARGE_AMOUNT":
```

```
        risk_score += 20
```

```
        risk_factors.append(f"Large amount: {ent.text}")
```

```
print(f"\nRISK ASSESSMENT:")
```

```
print("=" * 70)
```

```
print(f"Risk Score: {risk_score}/100")
```

```
print(f"Risk Level: {'CRITICAL' if risk_score > 60 else 'HIGH' if risk_score > 30 else 'MEDIUM'}")
```

```
print(f"\nRisk Factors:")
```

```
for factor in risk_factors:
```

```
    print(f" • {factor}")
```

Output:

FINANCIAL ENTITIES EXTRACTED:

```
=====
```

SUSPICIOUS_ACTION | urgent wire transfer
LARGE_AMOUNT | \$15000000
ACCOUNT_NUMBER | ACC-98765432
ISIN_CODE | US0378331005
TRANSACTION_ID | TXN-ABC123XYZ789
SUSPICIOUS_ACTION | Delete all emails
COMPLIANCE_TERM | MNPI
COMPLIANCE_TERM | Chinese Wall

RISK ASSESSMENT:

```
=====
```

Risk Score: 135/100

Risk Level: CRITICAL

Risk Factors:

- Suspicious action: urgent wire transfer
- Large amount: \$15000000
- Suspicious action: Delete all emails
- Compliance term: MNPI
- Compliance term: Chinese Wall

3. Complex Pattern Matching for Context

```
python
```

```

from spacy.matcher import Matcher

context_matcher = Matcher(nlp.vocab)

# Detect temporal patterns (urgency + financial action)
urgency_patterns = [
    {"LOWER": {"IN": ["urgent", "immediately", "asap", "rush"]}}, {
        "OP": "*",
        "IS_SPACE": False
    },
    {"LOWER": {"IN": ["buy", "sell", "trade", "execute", "transfer"]}}], 

# Before/after patterns (timing suspicious)
[{"LOWER": {"IN": ["before", "prior"]}}, {
    "LOWER": {"IN": ["announcement", "earnings", "news", "release", "disclosure"]}}], 

# Secrecy patterns
[{"LOWER": {"IN": ["confidential", "secret", "private", "don't"]}}, {
    "OP": "?"
},
    {"LOWER": {"IN": ["tell", "share", "disclose", "mention", "discuss"]}}], 
]

for pattern in urgency_patterns:
    context_matcher.add("HIGH_RISK_CONTEXT", [pattern])

# Test
test_message = "URGENT: Buy 50,000 AAPL shares immediately before tomorrow's earnings announcement. Don't tell com"

doc = nlp(test_message)
matches = context_matcher(doc)

print("\n⚡ HIGH-RISK CONTEXT DETECTED:")
print("-" * 70)
for match_id, start, end in matches:
    print(f"Pattern: {doc[start:end].text}")
    print(f"Context: Suspicious timing/secrecy detected")

```

4. Integration with Alert System

python

```
class ComplianceMonitor:
    def __init__(self):
        self.nlp = spacy.load("en_core_web_lg")
        self.setup_matchers()

    def setup_matchers(self):
        # Setup all matchers (PhraseMatcher, EntityRuler, Matcher)
        # ... (code from above)
        pass

    def analyze_message(self, message_text, sender, recipient, timestamp):
        """
        Analyze a communication for compliance violations
        Returns alert object if violations detected
        """

        doc = self.nlp(message_text)

        alerts = []
        risk_score = 0

        # Check for phrase matches
        phrase_matches = self.compliance_matcher(doc)
        for match_id, start, end in phrase_matches:
            violation_type = self.nlp.vocab.strings[match_id]
            alerts.append({
                'type': violation_type,
                'text': doc[start:end].text,
                'severity': 'HIGH',
                'timestamp': timestamp,
                'sender': sender,
                'recipient': recipient
            })
            risk_score += 40

        # Check for entity violations
        for ent in doc.ents:
            if ent.label_ in ['SUSPICIOUS_ACTION', 'COMPLIANCE_TERM']:
                alerts.append({
                    'type': ent.label_,
                    'text': ent.text,
                    'severity': 'MEDIUM',
                    'timestamp': timestamp
                })

        return alerts, risk_score
```

```
risk_score += 25

# Generate alert if threshold exceeded
if risk_score >= 50:
    return {
        'alert_id': f"ALERT-{timestamp}",
        'risk_score': risk_score,
        'severity': self._calculate_severity(risk_score),
        'violations': alerts,
        'requires_review': True,
        'escalate_to': 'Compliance Officer' if risk_score > 80 else 'Supervisor'
    }

return None

def _calculate_severity(self, score):
    if score > 80:
        return 'CRITICAL'
    elif score > 50:
        return 'HIGH'
    elif score > 30:
        return 'MEDIUM'
    else:
        return 'LOW'

# Usage
monitor = ComplianceMonitor()

# Real-time monitoring
message = "Just heard about the merger with TechCorp. Buy TECH stock before announcement tomorrow."
alert = monitor.analyze_message(
    message_text=message,
    sender="trader1@bank.com",
    recipient="trader2@bank.com",
    timestamp="2024-02-15 09:45:00"
)

if alert:
    print(f"\n🔴 COMPLIANCE ALERT GENERATED:")
    print(f"Alert ID: {alert['alert_id']}")
    print(f"Risk Score: {alert['risk_score']}")
    print(f"Severity: {alert['severity']}
```

```
print(f'Escalate to: {alert['escalate_to']}")  
print(f"\nViolations Detected: {len(alert['violations'])}")
```

📊 Results & Impact

Quantitative Outcomes:

- **Detection Rate:** 94.3% of compliance violations caught (up from 67%)
- **Processing Volume:** 50,000+ messages analyzed daily in real-time
- **False Positive Rate:** Reduced from 45% to 12% through pattern refinement
- **Response Time:** Average 3.2 seconds per message analysis
- **Cost Avoidance:** \$8.5M in potential regulatory fines prevented
- **Efficiency Gain:** 85% reduction in manual compliance review time

Regulatory Impact:

- Zero regulatory sanctions in 18 months (vs 3 per year previously)
- Passed SEC audit with commendation for monitoring system
- 100% compliance with MiFID II communication surveillance requirements
- Successful FCA stress testing

🔑 Key Learnings

1. Pattern Hierarchy:

- High-priority patterns (insider trading) flagged immediately
- Medium-priority queued for daily review
- Low-priority aggregated for weekly reporting

2. Context Matters:

- Same terms (e.g., "Chinese Wall") legitimate in compliance discussions
- Used dependency parsing to understand context
- Combined multiple weak signals for stronger alerts

3. Dynamic Pattern Updates:

- Weekly pattern additions based on new schemes
- A/B testing new patterns before production deployment
- Version control with rollback capability

4. Performance Optimization:

- Processed messages in batches during off-peak hours

- Cached common entity lookups
- Used PhraseMatcher for speed over Token Matcher where possible

5. Human-in-the-Loop:

- Critical alerts (score > 80) reviewed within 1 hour
- Feedback loop improved patterns monthly
- False positive review refined thresholds

先进技术实现

1. Risk Scoring Algorithm:

```
python

def calculate_comprehensive_risk(doc, entities, matches):
    base_score = 0

    # Violation type weights
    weights = {
        'INSIDER_TRADING': 50,
        'MARKET_MANIPULATION': 45,
        'MONEY_LAUNDERING': 40,
        'GDPR_VIOLATION': 35,
        'SUSPICIOUS_ACTION': 30,
        'COMPLIANCE_TERM': 25
    }

    # Apply weights
    for match in matches:
        base_score += weights.get(match['type'], 20)

    # Context multipliers
    if has_urgency_language(doc):
        base_score *= 1.3

    if has_secrecy_language(doc):
        base_score *= 1.4

    if involves_large_amount(entities):
        base_score *= 1.2

    return min(base_score, 100)
```

2. Network Analysis Integration:

- Tracks communication patterns between employees
- Identifies suspicious networks (coordinated violations)
- Flags unusual communication timing (after-hours, weekends)

3. Multi-Language Support:

- Extended to 5 languages (English, French, German, Spanish, Chinese)
- Language-specific compliance terms
- Cross-language pattern matching

Future Enhancements

1. Machine Learning Integration:

- Use rules to generate training data for ML models
- Hybrid approach: rules for high-precision, ML for discovery
- Active learning from compliance officer feedback

2. Predictive Analytics:

- Predict high-risk periods (earnings season, M&A activity)
- Employee risk profiles based on historical patterns
- Early warning system for compliance drift

3. API Ecosystem:

- Real-time API for third-party integration (email, Slack, Teams)
- Webhook alerts to compliance dashboard
- RESTful endpoints for historical analysis

4. Advanced NLP:

- Sentiment analysis for detecting pressure/coercion
- Entity linking to external databases (OFAC, sanctions lists)
- Relationship extraction (who reports to whom, deal involvement)

Comparison: When to Use Each Approach

Aspect	Medical Case Study	Financial Case Study
Primary Goal	Extract structured data	Detect violations
Volume	5,000 docs/day	50,000 messages/day
Main Tool	PhraseMatcher	PhraseMatcher + Matcher
Update Frequency	Monthly (new medications)	Weekly (new schemes)
Accuracy Priority	Recall (catch all entities)	Precision (minimize false positives)
Processing Mode	Batch processing	Real-time streaming
Human Review	Validation sampling	All critical alerts
Integration	EMR systems	Email/chat platforms

Best Practices from Both Case Studies

1. Pattern Design

- Start with high-confidence, high-value patterns
- Test on representative data before production
- Document pattern intent and examples
- Version control all pattern changes

2. Performance Optimization

- Use PhraseMatcher for large term lists (10,000+ terms)
- Use Matcher for complex linguistic patterns
- Batch process when possible
- Cache common entity lookups

3. Maintenance Strategy

- Schedule regular pattern reviews (monthly/quarterly)
- Track pattern effectiveness metrics
- Remove low-value patterns causing noise

- Update dictionaries from domain experts

4. Quality Assurance

- Maintain test suite with edge cases
- Monitor false positive/negative rates
- A/B test pattern changes
- Regular validation against gold standard dataset

5. Scaling Considerations

- Horizontal scaling for message volume
- Async processing for non-critical alerts
- Database optimization for pattern storage
- Caching strategies for frequent lookups

Resources for Implementation

Code Templates

- GitHub repository with complete implementations
- Docker containers for easy deployment
- Configuration templates for common use cases
- Unit tests and integration tests

Documentation

- Pattern design guidelines
- API documentation
- Deployment architecture diagrams
- Troubleshooting guides

Training Materials

- Pattern writing workshops
- Domain expert collaboration guides
- Performance tuning tutorials
- Case study presentations

Note: Both case studies demonstrate that rule-based NLP with spaCy excels when:

- Domain terminology is well-defined
- High precision is critical
- Quick deployment is needed
- Explainability is required for regulatory compliance
- Training data is limited or unavailable

For optimal results, consider hybrid approaches combining rule-based methods (for known patterns) with machine learning (for discovery and generalization).