

Structured Output - Advanced Handout

Machine Learning for Smarter Innovation

1 Structured Output - Advanced Handout

Target Audience: Data scientists and ML engineers **Duration:** 90 minutes reading **Level:** Advanced (mathematical foundations + production systems)

1.1 Mathematical Foundations

1.1.1 Language Models and Output Generation

Language models generate text by modeling the probability distribution over sequences. For a sequence of tokens w_1, w_2, \dots, w_n , an autoregressive language model factorizes the joint probability as:

$$P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$$

At each generation step, the model produces a probability distribution over the vocabulary \mathcal{V} . The next token is sampled or selected according to:

$$w_i \sim \text{softmax}\left(\frac{z_i}{\tau}\right)$$

where $z_i \in \mathbb{R}^{|\mathcal{V}|}$ are the logits and τ is the temperature parameter. Lower temperatures sharpen the distribution toward deterministic selection.

1.1.2 Information-Theoretic View of Structure

Structured output constrains the space of valid outputs from the full token sequence space to a strict subset. Let \mathcal{S} denote the set of valid structured outputs (e.g., well-formed JSON matching a schema). The entropy of the constrained distribution is:

$$H(\mathcal{S}) = - \sum_{s \in \mathcal{S}} P(s) \log P(s)$$

Constraining output structure reduces entropy, which can be quantified by the information gain:

$$\Delta I = H(\mathcal{V}^*) - H(\mathcal{S})$$

where \mathcal{V}^* is the set of all possible token sequences. This reduction in entropy translates to increased reliability and predictability of outputs.

1.1.3 Schema Theory and Type Systems

A JSON schema defines a formal grammar over structured documents. Schemas can be expressed in Backus-Naur Form (BNF):

```

<object> ::= "{" <members> "}" | "{}"
<members> ::= <pair> | <pair> "," <members>
<pair> ::= <string> ":" <value>
<value> ::= <object> | <array> | <string> | <number> | <boolean> | "null"
<array> ::= "[" <elements> "]" | "[]"
<elements> ::= <value> | <value> "," <elements>

```

Type constraints further restrict valid values. For a schema with property constraints, the validation function $V : \mathcal{D} \rightarrow \{0, 1\}$ is:

$$V(d) = \prod_{p \in \mathcal{P}} \mathbb{1}[\text{type}(d.p) \in \mathcal{T}_p \wedge \text{constraint}(d.p)]$$

where \mathcal{P} is the set of properties, \mathcal{T}_p is the allowed type set for property p , and $\text{constraint}(d.p)$ checks additional constraints (ranges, patterns, enumerations).

1.2 Constrained Decoding Theory

1.2.1 Grammar-Constrained Generation

Grammar-constrained decoding modifies the sampling process to only allow tokens that maintain grammatical validity. Given a context-free grammar $G = (N, \Sigma, P, S)$ where: - N is the set of non-terminal symbols - Σ is the terminal alphabet (vocabulary tokens) - P is the set of production rules - S is the start symbol

At each generation step, we compute the set of valid next tokens:

$$\mathcal{V}_{\text{valid}}(s) = \{v \in \Sigma : \exists w \in \Sigma^* \text{ such that } s \cdot v \cdot w \in L(G)\}$$

where s is the current prefix and $L(G)$ is the language defined by grammar G . The constrained distribution becomes:

$$P_{\text{constrained}}(w_i | w_{<i}) = \frac{P(w_i | w_{<i}) \cdot \mathbb{1}[w_i \in \mathcal{V}_{\text{valid}}(w_{<i})]}{\sum_{v \in \mathcal{V}_{\text{valid}}(w_{<i})} P(v | w_{<i})}$$

1.2.2 Finite State Machine Approach

JSON structure can be enforced using a finite state machine (FSM). The FSM $M = (Q, \Sigma, \delta, q_0, F)$ tracks the parsing state: - Q : states representing parser positions - Σ : input alphabet (tokens) - $\delta : Q \times \Sigma \rightarrow Q$: transition function - q_0 : initial state - $F \subseteq Q$: accepting states

For JSON generation, states include: OBJECT_START, KEY_EXPECTED, COLON_EXPECTED, VALUE_EXPECTED, COMMA_OR_END, ARRAY_START, etc.

The valid token mask at state q is:

$$\text{mask}(q) = \{v \in \Sigma : \delta(q, v) \neq \emptyset\}$$

1.2.3 Incremental Parsing Algorithm

Efficient constrained decoding requires incremental parsing that maintains a stack of parser states:

Algorithm: Incremental JSON Parser

```

Input: token sequence w_1, ..., w_n
Output: valid boolean, parser state

state_stack = [INITIAL]
for i = 1 to n:
    current_state = peek(state_stack)
    if w_i not in valid_tokens(current_state):
        return False, ERROR

    new_states = transition(current_state, w_i)
    update_stack(state_stack, new_states)

return is_accepting(state_stack), state_stack

```

The time complexity is $O(n)$ for fixed-depth JSON structures, with space complexity $O(d)$ where d is the maximum nesting depth.

1.3 Probabilistic Schema Matching

1.3.1 Maximum Likelihood Extraction

Given an input text x and a target schema S , structured extraction finds the most likely structured output:

$$\hat{y} = \underset{y \in \mathcal{Y}_S}{\operatorname{argmax}} P(y|x, S)$$

where \mathcal{Y}_S is the set of valid outputs conforming to schema S . Using Bayes' theorem:

$$P(y|x, S) \propto P(x|y, S) \cdot P(y|S)$$

The prior $P(y|S)$ encodes schema constraints, assigning zero probability to invalid structures.

1.3.2 Confidence Estimation

Extraction confidence can be estimated from the model's probability distribution. For a generated field value v_i , the confidence is:

$$c_i = \frac{1}{|v_i|} \sum_{j=1}^{|v_i|} \log P(t_{i,j} | t_{i,<j}, x)$$

where $t_{i,j}$ are the tokens comprising value v_i . The overall extraction confidence aggregates field confidences:

$$c_{\text{total}} = \frac{1}{n} \sum_{i=1}^n w_i \cdot c_i$$

with field weights w_i reflecting importance. Fields with low confidence should trigger validation or human review.

1.3.3 Calibration Theory

Confidence scores require calibration to reflect true reliability. A perfectly calibrated model satisfies:

$$P(\text{correct}|c) = c$$

Calibration error is measured using Expected Calibration Error (ECE):

$$\text{ECE} = \sum_{b=1}^B \frac{|B_b|}{n} |acc(B_b) - conf(B_b)|$$

where B_b are bins of predictions grouped by confidence, $acc(B_b)$ is accuracy within bin b , and $conf(B_b)$ is mean confidence.

Platt scaling calibrates raw scores using logistic regression:

$$c_{\text{calibrated}} = \sigma(a \cdot c_{\text{raw}} + b)$$

where parameters a, b are learned on a held-out calibration set.

1.4 Advanced Validation Theory

1.4.1 Multi-Layer Validation Architecture

Production systems require multiple validation layers with different error detection capabilities:

Layer 1: Syntactic Validation - Checks: well-formed JSON, valid types - Error rate: detects ~60% of issues - Cost: $O(n)$ parsing, negligible

Layer 2: Schema Validation - Checks: required fields, type constraints, patterns - Error rate: detects ~85% of issues - Cost: $O(n)$ traversal, minimal

Layer 3: Semantic Validation - Checks: value consistency, business rules - Error rate: detects ~95% of issues - Cost: $O(n)$ to $O(n^2)$ depending on rules

Layer 4: Statistical Validation - Checks: distributional anomalies, outliers - Error rate: detects ~99% of issues - Cost: requires historical data, $O(1)$ lookup

The cumulative detection probability after k layers:

$$P(\text{detect}) = 1 - \prod_{i=1}^k (1 - p_i)$$

where p_i is the detection probability of layer i .

1.4.2 Consistency Constraints

Cross-field consistency rules enforce logical relationships. Define a constraint function $C : \mathcal{Y} \rightarrow \{0, 1\}$:

$$C(y) = \bigwedge_{(i,j) \in \mathcal{R}} \phi_{i,j}(y.f_i, y.f_j)$$

where \mathcal{R} is the set of related field pairs and $\phi_{i,j}$ is the consistency predicate.

Example constraints: - **Range consistency**: if rating = 5, then satisfaction ≥ 4 - **Temporal consistency**: if end_date exists, then end_date \geq start_date - **Logical consistency**: if status = "completed", then completion_date is not null

1.5 Implementation: Production Validation Pipeline

```
"""
Production-grade structured output validation system.
Implements multi-layer validation with confidence scoring.
"""

from typing import Protocol, List, Dict, Any, Optional, Tuple
from dataclasses import dataclass, field
from abc import ABC, abstractmethod
from enum import Enum
import json
import re
import numpy as np
from pydantic import BaseModel, ValidationError, field_validator
from scipy import stats

class ValidationSeverity(Enum):
    """Severity levels for validation issues."""
    ERROR = "error" # Hard failure, reject output
    WARNING = "warning" # Soft failure, flag for review
    INFO = "info" # Informational, log only

@dataclass
class ValidationIssue:
    """Single validation issue with context."""
    severity: ValidationSeverity
    field: str
    message: str
    expected: Optional[Any] = None
    actual: Optional[Any] = None
    confidence: float = 1.0

@dataclass
class ValidationResult:
    """Complete validation result with aggregated metrics."""
    valid: bool
    confidence: float
    issues: List[ValidationIssue] = field(default_factory=list)
    metadata: Dict[str, Any] = field(default_factory=dict)

    @property
    def error_count(self) -> int:
        return sum(1 for i in self.issues if i.severity == ValidationSeverity.ERROR)

    @property
    def warning_count(self) -> int:
        return sum(1 for i in self.issues if i.severity == ValidationSeverity.WARNING)

class Validator(ABC):
    """Abstract base class for validators."""

    @abstractmethod
    def validate(self, data: Dict[str, Any]) -> ValidationResult:
```

```

    """Validate data and return result."""
    pass

    @property
    @abstractmethod
    def name(self) -> str:
        """Validator name for logging."""
        pass

class SchemaValidator(Validator):
    """Pydantic-based schema validation."""

    def __init__(self, schema_class: type):
        self.schema_class = schema_class

    @property
    def name(self) -> str:
        return f"SchemaValidator({self.schema_class.__name__})"

    def validate(self, data: Dict[str, Any]) -> ValidationResult:
        issues = []
        try:
            self.schema_class(**data)
            return ValidationResult(valid=True, confidence=1.0, issues=[])
        except ValidationError as e:
            for error in e.errors():
                issues.append(ValidationIssue(
                    severity=ValidationSeverity.ERROR,
                    field=".".join(str(loc) for loc in error["loc"]),
                    message=error["msg"],
                    actual=error.get("input")
                ))
            return ValidationResult(valid=False, confidence=0.0, issues=issues
)

class BusinessRuleValidator(Validator):
    """Configurable business rule validation."""

    def __init__(self, rules: List[Dict[str, Any]]):
        self.rules = rules

    @property
    def name(self) -> str:
        return f"BusinessRuleValidator({len(self.rules)} rules)"

    def validate(self, data: Dict[str, Any]) -> ValidationResult:
        issues = []

        for rule in self.rules:
            if not self._evaluate_rule(data, rule):
                issues.append(ValidationIssue(
                    severity=rule.get("severity", ValidationSeverity.WARNING),
                    field=rule.get("field", ""),
                    message=rule["message"],
                    expected=rule.get("expected"),
                    actual=self._get_nested(data, rule.get("field", ""))
                ))

        has_errors = any(i.severity == ValidationSeverity.ERROR for i in
issues)
        confidence = 1.0 - (len(issues) * 0.1) # Degrade confidence per issue

```

```

    return ValidationResult(
        valid=not has_errors,
        confidence=max(0.0, confidence),
        issues=issues
    )

def _evaluate_rule(self, data: Dict, rule: Dict) -> bool:
    """Evaluate a single rule against data."""
    condition = rule.get("condition", lambda d: True)
    return condition(data)

def _get_nested(self, data: Dict, path: str) -> Any:
    """Get nested value by dot-separated path."""
    if not path:
        return None
    parts = path.split(".")
    current = data
    for part in parts:
        if isinstance(current, dict):
            current = current.get(part)
        else:
            return None
    return current

class StatisticalValidator(Validator):
    """Statistical anomaly detection for numeric fields."""

    def __init__(self, field_statistics: Dict[str, Dict[str, float]]):
        """
        Initialize with historical statistics per field.

        Args:
            field_statistics: Dict mapping field names to stats
                {"field": {"mean": x, "std": y, "min": a, "max":
b}}
        """
        self.stats = field_statistics
        self.z_threshold = 3.0 # Standard deviations for anomaly

    @property
    def name(self) -> str:
        return f"StatisticalValidator({len(self.stats)} fields)"

    def validate(self, data: Dict[str, Any]) -> ValidationResult:
        issues = []
        confidences = []

        for field_name, field_stats in self.stats.items():
            value = self._get_nested(data, field_name)
            if value is None or not isinstance(value, (int, float)):
                continue

            # Z-score calculation
            mean = field_stats.get("mean", 0)
            std = field_stats.get("std", 1)
            z_score = abs((value - mean) / std) if std > 0 else 0

            # Range check
            min_val = field_stats.get("min", float("-inf"))
            max_val = field_stats.get("max", float("inf"))

```

```

        if value < min_val or value > max_val:
            issues.append(ValidationIssue(
                severity=ValidationSeverity.ERROR,
                field=field_name,
                message=f"Value {value} outside valid range [{min_val}, {
max_val}]",
                expected=f"[{min_val}, {max_val}]",
                actual=value
            ))
        elif z_score > self.z_threshold:
            issues.append(ValidationIssue(
                severity=ValidationSeverity.WARNING,
                field=field_name,
                message=f"Value {value} is {z_score:.1f} std devs from
mean {mean:.1f}",
                confidence=1 - stats.norm.cdf(z_score)
            ))

        # Confidence based on z-score
        field_confidence = 1 - min(z_score / (2 * self.z_threshold), 1.0)
        confidences.append(field_confidence)

    overall_confidence = np.mean(confidences) if confidences else 1.0
    has_errors = any(i.severity == ValidationSeverity.ERROR for i in
issues)

    return ValidationResult(
        valid=not has_errors,
        confidence=overall_confidence,
        issues=issues
    )

def _get_nested(self, data: Dict, path: str) -> Any:
    parts = path.split(".")
    current = data
    for part in parts:
        if isinstance(current, dict):
            current = current.get(part)
        else:
            return None
    return current

class ValidationPipeline:
    """
    Multi-stage validation pipeline with short-circuit on errors.
    """

    def __init__(self, validators: List[Validator], fail_fast: bool = True):
        self.validators = validators
        self.fail_fast = fail_fast

    def validate(self, data: Dict[str, Any]) -> ValidationResult:
        """Run all validators and aggregate results."""
        all_issues = []
        min_confidence = 1.0
        metadata = {"validators_run": [], "validators_passed": []}

        for validator in self.validators:
            result = validator.validate(data)
            metadata["validators_run"].append(validator.name)

            all_issues.extend(result.issues)

```

```

        min_confidence = min(min_confidence, result.confidence)

    if result.valid:
        metadata["validators_passed"].append validator.name
    elif self.fail_fast:
        # Stop on first hard failure
        return ValidationResult(
            valid=False,
            confidence=0.0,
            issues=all_issues,
            metadata=metadata
        )

    has_errors = any(i.severity == ValidationSeverity.ERROR for i in
all_issues)

    return ValidationResult(
        valid=not has_errors,
        confidence=min_confidence,
        issues=all_issues,
        metadata=metadata
    )

```

1.6 Implementation: Constrained Generation Engine

```

"""
Grammar-constrained JSON generation using finite state machines.
Ensures 100% syntactically valid output.
"""

from typing import Dict, Set, List, Optional, Callable
from dataclasses import dataclass
from enum import Enum, auto
import json

class JSONState(Enum):
    """Parser states for JSON generation."""
    START = auto()
    OBJECT_OPEN = auto()
    OBJECT_KEY = auto()
    OBJECT_COLON = auto()
    OBJECT_VALUE = auto()
    OBJECT_COMMA = auto()
    ARRAY_OPEN = auto()
    ARRAY_VALUE = auto()
    ARRAY_COMMA = auto()
    STRING_CONTENT = auto()
    NUMBER = auto()
    COMPLETE = auto()

@dataclass
class SchemaField:
    """Schema field definition."""
    name: str
    type: str # "string", "number", "boolean", "object", "array"
    required: bool = True
    enum: Optional[List[str]] = None

```

```

minimum: Optional[float] = None
maximum: Optional[float] = None
properties: Optional[Dict] = None # For nested objects

class JSONConstrainedDecoder:
    """
    Finite state machine for constrained JSON generation.
    Tracks valid tokens at each generation step.
    """

    def __init__(self, schema: Dict[str, SchemaField]):
        self.schema = schema
        self.state_stack: List[JSONState] = [JSONState.START]
        self.field_stack: List[str] = []
        self.generated: List[str] = []
        self.current_field_idx = 0
        self.field_order = list(schema.keys())

    def get_valid_tokens(self) -> Set[str]:
        """Return set of valid next tokens given current state."""
        state = self.state_stack[-1]

        if state == JSONState.START:
            return {"{"}

        elif state == JSONState.OBJECT_OPEN:
            if self.current_field_idx < len(self.field_order):
                field_name = self.field_order[self.current_field_idx]
                return {f'"{field_name}"'}
            return {"}"}

        elif state == JSONState.OBJECT_KEY:
            return {":"}

        elif state == JSONState.OBJECT_COLON:
            return self._valid_value_tokens()

        elif state == JSONState.OBJECT_VALUE:
            if self.current_field_idx < len(self.field_order):
                return {","}
            return {"}"}

        elif state == JSONState.OBJECT_COMMA:
            field_name = self.field_order[self.current_field_idx]
            return {f'"{field_name}"'}

        elif state == JSONState.ARRAY_OPEN:
            return self._valid_value_tokens() | {"["}

        elif state == JSONState.ARRAY_VALUE:
            return {"", ", ", "]" }

        elif state == JSONState.COMPLETE:
            return set()

        return set()

    def _valid_value_tokens(self) -> Set[str]:
        """Get valid tokens for value position based on field schema."""
        if not self.field_stack:
            field_name = self.field_order[self.current_field_idx - 1]
        else:

```

```

        field_name = self.field_stack[-1]

    field = self.schema.get(field_name)
    if not field:
        return {"null"}

    if field.type == "string":
        if field.enum:
            return {f'"{v}"' for v in field.enum}
        return {"<STRING>"} # Placeholder for string content

    elif field.type == "number":
        return {"<NUMBER>"} # Placeholder for numeric content

    elif field.type == "boolean":
        return {"true", "false"}

    elif field.type == "object":
        return {"{"}

    elif field.type == "array":
        return {"["}

    return {"null"}

def advance(self, token: str) -> bool:
    """
    Advance state machine with given token.
    Returns True if token was valid, False otherwise.
    """
    valid_tokens = self.get_valid_tokens()

    # Handle placeholder tokens
    if "<STRING>" in valid_tokens and token.startswith('"') and token.
endswith('"'):
        pass # Accept any string
    elif "<NUMBER>" in valid_tokens and self._is_number(token):
        pass # Accept any number
    elif token not in valid_tokens:
        return False

    self.generated.append(token)
    self._update_state(token)
    return True

def _update_state(self, token: str):
    """Update state stack based on consumed token."""
    state = self.state_stack[-1]

    if token == "{":
        self.state_stack.append(JSONState.OBJECT_OPEN)
    elif token == "}":
        self.state_stack.pop()
        if self.state_stack[-1] == JSONState.OBJECT_COLON:
            self.state_stack[-1] = JSONState.OBJECT_VALUE
        elif self.state_stack[-1] == JSONState.ARRAY_VALUE:
            pass
        elif len(self.state_stack) == 1:
            self.state_stack[-1] = JSONState.COMPLETE
    elif token == "[":
        self.state_stack.append(JSONState.ARRAY_OPEN)
    elif token == "]":
        self.state_stack.pop()

```

```

        if self.state_stack[-1] == JSONState.OBJECT_COLON:
            self.state_stack[-1] = JSONState.OBJECT_VALUE
    elif token == ":":
        self.state_stack[-1] = JSONState.OBJECT_COLON
    elif token == ",":
        if self.state_stack[-1] == JSONState.OBJECT_VALUE:
            self.current_field_idx += 1
            self.state_stack[-1] = JSONState.OBJECT_COMMA
        elif self.state_stack[-1] == JSONState.ARRAY_VALUE:
            self.state_stack[-1] = JSONState.ARRAY_OPEN
    elif token.startswith('"') and state == JSONState.OBJECT_OPEN:
        self.current_field_idx += 1
        self.state_stack[-1] = JSONState.OBJECT_KEY
    elif token.startswith('"') and state == JSONState.OBJECT_COMMA:
        self.state_stack[-1] = JSONState.OBJECT_KEY
    elif state == JSONState.OBJECT_COLON:
        self.state_stack[-1] = JSONState.OBJECT_VALUE
    elif state == JSONState.ARRAY_OPEN:
        self.state_stack[-1] = JSONState.ARRAY_VALUE

def _is_number(self, token: str) -> bool:
    try:
        float(token)
        return True
    except ValueError:
        return False

def is_complete(self) -> bool:
    """Check if generation is complete and valid."""
    return self.state_stack[-1] == JSONState.COMPLETE

def get_generated_json(self) -> str:
    """Return the generated JSON string."""
    return "".join(self.generated)

class ConstrainedGenerator:
    """
    High-level constrained generation using LLM with FSM guidance.
    """

    def __init__(self, llm_client, schema: Dict[str, SchemaField]):
        self.llm = llm_client
        self.schema = schema

    def generate(self, prompt: str, max_tokens: int = 500) -> Dict:
        """
        Generate structured output with grammar constraints.
        """
        decoder = JSONConstrainedDecoder(self.schema)

        full_prompt = self._build_prompt(prompt)

        for _ in range(max_tokens):
            valid_tokens = decoder.get_valid_tokens()
            if not valid_tokens:
                break

            # Get next token from LLM with constrained sampling
            next_token = self._sample_constrained(full_prompt, valid_tokens)

            if not decoder.advance(next_token):
                # Fallback: pick first valid token

```

```

        next_token = list(valid_tokens)[0]
        decoder.advance(next_token)

    full_prompt += next_token

    if decoder.is_complete():
        break

    return json.loads(decoder.get_generated_json())

def _build_prompt(self, user_prompt: str) -> str:
    """Build prompt with schema instructions."""
    schema_desc = self._schema_to_description()
    return f"""Extract information into JSON format.

Schema:
{schema_desc}

Input: {user_prompt}

Output JSON:
"""

def _schema_to_description(self) -> str:
    """Convert schema to human-readable description."""
    lines = []
    for name, field in self.schema.items():
        constraint = ""
        if field.enum:
            constraint = f" (one of: {', '.join(field.enum)})"
        elif field.minimum is not None or field.maximum is not None:
            constraint = f" (range: {field.minimum or '-inf'} to {field.
maximum or 'inf'})"
        lines.append(f"- {name}: {field.type}{constraint}")
    return "\n".join(lines)

def _sample_constrained(self, prompt: str, valid_tokens: Set[str]) -> str:
    """Sample from LLM with token constraints."""
    # Implementation depends on LLM API
    # This is a simplified version
    response = self.llm.generate(
        prompt=prompt,
        max_tokens=10,
        logit_bias=self._tokens_to_bias(valid_tokens)
    )
    return response.strip()

def _tokens_to_bias(self, valid_tokens: Set[str]) -> Dict[int, float]:
    """Convert valid tokens to logit bias dictionary."""
    # Map token strings to token IDs and apply positive bias
    # Negative infinity bias for invalid tokens
    return {} # Placeholder - implement based on tokenizer

```

1.7 Implementation: Multi-Model Extraction System

```

"""
Production multi-model extraction with fallback and ensemble.
Achieves 99%+ reliability through redundancy.
"""

```

```

from typing import Dict, List, Tuple, Optional, Callable
from dataclasses import dataclass
from enum import Enum
import asyncio
import time
import logging
from abc import ABC, abstractmethod

logger = logging.getLogger(__name__)

class ModelTier(Enum):
    """Model tiers ordered by capability and cost."""
    PREMIUM = "premium" # GPT-4, Claude-3-Opus
    STANDARD = "standard" # GPT-3.5, Claude-3-Sonnet
    FALLBACK = "fallback" # Local models, rule-based
    ENSEMBLE = "ensemble" # Combine multiple models

@dataclass
class ExtractionResult:
    """Result from a single extraction attempt."""
    data: Optional[Dict]
    confidence: float
    model_tier: ModelTier
    latency_ms: float
    error: Optional[str] = None

class ModelClient(ABC):
    """Abstract base for model clients."""

    @abstractmethod
    async def extract(self, prompt: str, schema: Dict) -> ExtractionResult:
        pass

    @property
    @abstractmethod
    def tier(self) -> ModelTier:
        pass

class OpenAIClient(ModelClient):
    """OpenAI API client."""

    def __init__(self, model: str = "gpt-4", api_key: str = None):
        self.model = model
        self.api_key = api_key
        self._tier = ModelTier.PREMIUM if "gpt-4" in model else ModelTier.STANDARD

    @property
    def tier(self) -> ModelTier:
        return self._tier

    async def extract(self, prompt: str, schema: Dict) -> ExtractionResult:
        start = time.time()
        try:
            # Actual API call would go here
            response = await self._call_api(prompt, schema)
            latency = (time.time() - start) * 1000

```

```

        return ExtractionResult(
            data=response,
            confidence=self._estimate_confidence(response),
            model_tier=self.tier,
            latency_ms=latency
        )
    except Exception as e:
        return ExtractionResult(
            data=None,
            confidence=0.0,
            model_tier=self.tier,
            latency_ms=(time.time() - start) * 1000,
            error=str(e)
        )

    async def _call_api(self, prompt: str, schema: Dict) -> Dict:
        # Placeholder for actual API call
        await asyncio.sleep(0.1)
        return {}

    def _estimate_confidence(self, response: Dict) -> float:
        # Estimate confidence based on response characteristics
        return 0.9

class MultiModelExtractor:
    """
    Orchestrates multiple models with fallback and ensemble strategies.
    """

    def __init__(self, clients: List[ModelClient], validation_pipeline):
        self.clients = sorted(clients, key=lambda c: c.tier.value)
        self.validation = validation_pipeline
        self.metrics = ExtractionMetrics()

    async def extract(
        self,
        text: str,
        schema: Dict,
        strategy: str = "fallback"
    ) -> Tuple[Dict, ExtractionResult]:
        """
        Extract structured data using specified strategy.

        Args:
            text: Input text to extract from
            schema: Target schema
            strategy: "fallback" (try in order) or "ensemble" (combine results)
        """

        Returns:
            Tuple of (extracted data, result metadata)
        """
        if strategy == "ensemble":
            return await self._ensemble_extract(text, schema)
        else:
            return await self._fallback_extract(text, schema)

    async def _fallback_extract(
        self,
        text: str,
        schema: Dict
    ) -> Tuple[Dict, ExtractionResult]:

```

```

"""Try each model until one succeeds validation."""
prompt = self._build_extraction_prompt(text, schema)

for client in self.clients:
    result = await client.extract(prompt, schema)
    self.metrics.record_attempt(client.tier, result)

    if result.error:
        logger.warning(f"{client.tier} failed: {result.error}")
        continue

    # Validate result
    validation = self.validation.validate(result.data)
    if validation.valid and validation.confidence > 0.7:
        logger.info(f"Success with {client.tier}")
        return result.data, result

    raise ExtractionError("All models failed")

async def _ensemble_extract(
    self,
    text: str,
    schema: Dict
) -> Tuple[Dict, ExtractionResult]:
    """Run multiple models in parallel and combine results."""
    prompt = self._build_extraction_prompt(text, schema)

    # Run all models concurrently
    tasks = [client.extract(prompt, schema) for client in self.clients]
    results = await asyncio.gather(*tasks, return_exceptions=True)

    # Filter successful results
    valid_results = []
    for result in results:
        if isinstance(result, Exception):
            continue
        if result.error or result.data is None:
            continue

        validation = self.validation.validate(result.data)
        if validation.valid:
            valid_results.append((result, validation.confidence))

    if not valid_results:
        raise ExtractionError("No valid results from ensemble")

    # Combine results using confidence-weighted voting
    combined = self._weighted_combine(valid_results, schema)

    ensemble_result = ExtractionResult(
        data=combined,
        confidence=np.mean([v[1] for v in valid_results]),
        model_tier=ModelTier.ENSEMBLE,
        latency_ms=max(r.latency_ms for r, _ in valid_results)
    )

    return combined, ensemble_result

def _weighted_combine(
    self,
    results: List[Tuple[ExtractionResult, float]],
    schema: Dict
) -> Dict:

```

```

    """Combine multiple extraction results using weighted voting."""
    combined = {}

    for field in schema.keys():
        values_with_weights = []
        for result, confidence in results:
            if field in result.data:
                values_with_weights.append((result.data[field], confidence
))

        if not values_with_weights:
            continue

        # For numeric fields, use weighted average
        if isinstance(values_with_weights[0][0], (int, float)):
            total_weight = sum(w for _, w in values_with_weights)
            combined[field] = sum(v * w for v, w in values_with_weights) /
total_weight
        else:
            # For categorical, use highest-confidence value
            combined[field] = max(values_with_weights, key=lambda x: x[1])
[0]

    return combined

def _build_extraction_prompt(self, text: str, schema: Dict) -> str:
    return f"Extract the following fields from the text:\n{text}"

class ExtractionError(Exception):
    """Raised when extraction fails across all models."""
    pass

@dataclass
class ExtractionMetrics:
    """Track extraction performance metrics."""
    attempts_by_tier: Dict[ModelTier, int] = None
    successes_by_tier: Dict[ModelTier, int] = None
    latencies_by_tier: Dict[ModelTier, List[float]] = None

    def __post_init__(self):
        self.attempts_by_tier = {t: 0 for t in ModelTier}
        self.successes_by_tier = {t: 0 for t in ModelTier}
        self.latencies_by_tier = {t: [] for t in ModelTier}

    def record_attempt(self, tier: ModelTier, result: ExtractionResult):
        self.attempts_by_tier[tier] += 1
        self.latencies_by_tier[tier].append(result.latency_ms)
        if result.data and not result.error:
            self.successes_by_tier[tier] += 1

    def get_success_rate(self, tier: ModelTier) -> float:
        if self.attempts_by_tier[tier] == 0:
            return 0.0
        return self.successes_by_tier[tier] / self.attempts_by_tier[tier]

    def get_p95_latency(self, tier: ModelTier) -> float:
        latencies = self.latencies_by_tier[tier]
        if not latencies:
            return 0.0
        return np.percentile(latencies, 95)

```

1.8 Implementation: Production Monitoring System

```

"""
Comprehensive monitoring and alerting for structured output systems.
Tracks reliability, latency, cost, and drift metrics.
"""

from typing import Dict, List, Callable, Optional
from dataclasses import dataclass, field
from datetime import datetime, timedelta
from collections import deque
import threading
import time

@dataclass
class MetricSample:
    """Single metric observation."""
    timestamp: datetime
    value: float
    tags: Dict[str, str] = field(default_factory=dict)

class MetricStore:
    """
    Time-series metric storage with configurable retention.
    """

    def __init__(self, retention_hours: int = 24):
        self.retention = timedelta(hours=retention_hours)
        self.metrics: Dict[str, deque] = {}
        self._lock = threading.Lock()

    def record(self, name: str, value: float, tags: Dict[str, str] = None):
        """Record a metric value."""
        sample = MetricSample(
            timestamp=datetime.now(),
            value=value,
            tags=tags or {}
        )

        with self._lock:
            if name not in self.metrics:
                self.metrics[name] = deque()
            self.metrics[name].append(sample)
            self._cleanup(name)

    def _cleanup(self, name: str):
        """Remove expired samples."""
        cutoff = datetime.now() - self.retention
        while self.metrics[name] and self.metrics[name][0].timestamp < cutoff:
            self.metrics[name].popleft()

    def get_recent(
        self,
        name: str,
        duration: timedelta = None,
        tags: Dict[str, str] = None
    ) -> List[MetricSample]:
        """Get recent samples, optionally filtered by tags."""

```

```

    duration = duration or timedelta(hours=1)
    cutoff = datetime.now() - duration

    with self._lock:
        samples = self.metrics.get(name, [])
        result = [
            s for s in samples
            if s.timestamp >= cutoff
            and (tags is None or all(s.tags.get(k) == v for k, v in tags.
items()))
        ]

    return result

def aggregate(
    self,
    name: str,
    aggregation: str = "mean",
    duration: timedelta = None
) -> float:
    """Compute aggregate statistic over recent samples."""
    samples = self.get_recent(name, duration)
    if not samples:
        return 0.0

    values = [s.value for s in samples]

    if aggregation == "mean":
        return sum(values) / len(values)
    elif aggregation == "sum":
        return sum(values)
    elif aggregation == "min":
        return min(values)
    elif aggregation == "max":
        return max(values)
    elif aggregation == "p50":
        return np.percentile(values, 50)
    elif aggregation == "p95":
        return np.percentile(values, 95)
    elif aggregation == "p99":
        return np.percentile(values, 99)
    elif aggregation == "count":
        return len(values)

    raise ValueError(f"Unknown aggregation: {aggregation}")

@dataclass
class AlertConfig:
    """Alert configuration."""
    name: str
    metric: str
    condition: Callable[[float], bool]
    severity: str # "critical", "warning", "info"
    message_template: str
    cooldown_minutes: int = 5

class AlertManager:
    """
    Alert evaluation and notification system.
    """

```

```

def __init__(self, metric_store: MetricStore):
    self.store = metric_store
    self.alerts: List[AlertConfig] = []
    self.last_fired: Dict[str, datetime] = {}
    self.handlers: Dict[str, List[Callable]] = {
        "critical": [],
        "warning": [],
        "info": []
    }

def register_alert(self, config: AlertConfig):
    """Register an alert configuration."""
    self.alerts.append(config)

def register_handler(self, severity: str, handler: Callable):
    """Register notification handler for severity level."""
    self.handlers[severity].append(handler)

def evaluate(self):
    """Evaluate all alerts against current metrics."""
    for alert in self.alerts:
        # Check cooldown
        if alert.name in self.last_fired:
            elapsed = datetime.now() - self.last_fired[alert.name]
            if elapsed < timedelta(minutes=alert.cooldown_minutes):
                continue

        # Get metric value
        value = self.store.aggregate(alert.metric, "mean", timedelta(
minutes=5))

        # Check condition
        if alert.condition(value):
            self._fire_alert(alert, value)

def _fire_alert(self, alert: AlertConfig, value: float):
    """Fire alert and notify handlers."""
    message = alert.message_template.format(value=value)
    self.last_fired[alert.name] = datetime.now()

    for handler in self.handlers[alert.severity]:
        try:
            handler(alert.name, message, value)
        except Exception as e:
            logger.error(f"Alert handler failed: {e}")

class StructuredOutputMonitor:
    """
    High-level monitoring interface for structured output systems.
    """

    def __init__(self):
        self.store = MetricStore(retention_hours=48)
        self.alerts = AlertManager(self.store)
        self._setup_default_alerts()

    def _setup_default_alerts(self):
        """Configure standard alerts."""
        self.alerts.register_alert(AlertConfig(
            name="low_success_rate",
            metric="extraction.success",
            condition=lambda v: v < 0.95,

```

```

        severity="critical",
        message_template="Extraction success rate dropped to {value:.1%}"
    ))

    self.alerts.register_alert(AlertConfig(
        name="high_latency",
        metric="extraction.latency_p95",
        condition=lambda v: v > 5000,
        severity="warning",
        message_template="P95 latency is {value:.0f}ms"
    ))

    self.alerts.register_alert(AlertConfig(
        name="validation_failures",
        metric="validation.failure_rate",
        condition=lambda v: v > 0.1,
        severity="warning",
        message_template="Validation failure rate is {value:.1%}"
    ))

def record_extraction(
    self,
    success: bool,
    latency_ms: float,
    model: str,
    confidence: float
):
    """Record extraction attempt metrics."""
    tags = {"model": model}

    self.store.record("extraction.success", 1.0 if success else 0.0, tags)
    self.store.record("extraction.latency", latency_ms, tags)
    self.store.record("extraction.confidence", confidence, tags)

def record_validation(self, passed: bool, error_count: int, warning_count:
int):
    """Record validation metrics."""
    self.store.record("validation.success", 1.0 if passed else 0.0)
    self.store.record("validation.errors", error_count)
    self.store.record("validation.warnings", warning_count)
    self.store.record("validation.failure_rate", 0.0 if passed else 1.0)

def record_cost(self, amount: float, model: str):
    """Record cost metrics."""
    self.store.record("cost.total", amount, {"model": model})

def get_dashboard_metrics(self) -> Dict:
    """Get metrics for dashboard display."""
    return {
        "success_rate": self.store.aggregate("extraction.success", "mean")
    ,
        "p50_latency": self.store.aggregate("extraction.latency", "p50"),
        "p95_latency": self.store.aggregate("extraction.latency", "p95"),
        "validation_rate": self.store.aggregate("validation.success", "
mean"),
        "total_cost": self.store.aggregate("cost.total", "sum"),
        "request_count": self.store.aggregate("extraction.success", "count
"),
        "avg_confidence": self.store.aggregate("extraction.confidence", "
mean")
    }

```

1.9 Common Parameters

Component	Parameter	Typical Range	Notes
Temperature	value	0.0-0.3	Lower = more deterministic
Max tokens	output limit	500-2000	Based on schema complexity
Retry count	attempts	2-5	Before fallback
Circuit breaker	failure threshold	3-10	Consecutive failures
Circuit breaker	timeout	30-300 sec	Recovery period
Cache TTL	duration	1-24 hours	Based on data freshness needs
Batch size	items	5-50	Balance throughput vs memory
Confidence threshold	minimum	0.7-0.9	For auto-accept
Human review threshold	range	0.5-0.7	Flag for manual check
Alert cooldown	minutes	5-60	Prevent alert fatigue

1.10 Practice Problems

- Constrained Decoder:** Implement a JSON decoder that uses beam search with grammar constraints. Track the valid token set at each step and compare generation quality vs unconstrained decoding.
 - Calibration Analysis:** Build a calibration curve for your extraction system. Collect 500+ samples, bin by confidence, and compute ECE. Implement Platt scaling and measure improvement.
 - Ensemble Optimization:** Compare fallback vs ensemble strategies across different input types. Measure when ensemble overhead is justified by accuracy gains.
 - Drift Detection:** Implement statistical drift detection for extraction outputs. Use Kolmogorov-Smirnov tests to compare output distributions over time windows.
 - Cost-Quality Tradeoff:** Build a router that selects models based on estimated input complexity. Optimize the complexity estimator to minimize cost while maintaining target accuracy.
-

1.11 References

- Willard, B. T., & Louf, R. (2023). "Efficient Guided Generation for Large Language Models." arXiv:2307.09702.
 - Geng, S., et al. (2023). "Grammar-Constrained Decoding for Structured NLP Tasks." ACL 2023.
 - Pydantic Documentation. "Data Validation Using Python Type Hints." <https://docs.pydantic.dev/>
 - OpenAI. (2024). "Function Calling and Structured Outputs." API Documentation.
 - Anthropic. (2024). "Tool Use and Structured Responses." Claude Documentation.
 - Guo, C., et al. (2017). "On Calibration of Modern Neural Networks." ICML 2017.
 - Nori, H., et al. (2023). "Capabilities of GPT-4 on Medical Challenge Problems." arXiv:2303.13375.
-

Structured output transforms probabilistic language models into deterministic data systems. The combination of grammar constraints, multi-layer validation, and production monitoring achieves the reliability required for enterprise applications.