# Structured Output
# & Reliable AI Systems

## From Prototype to Production

## Today's Journey

**Part 1: Foundation**
- The reliability challenge
- Prototype vs production gap
- Why structure matters
- Production requirements

**Part 2: Techniques**
- JSON schema fundamentals
- Prompt engineering patterns
- Function calling mechanics
- Validation strategies

**Part 3: Implementation**
- OpenAI function calling
- Pydantic validation
- Error handling
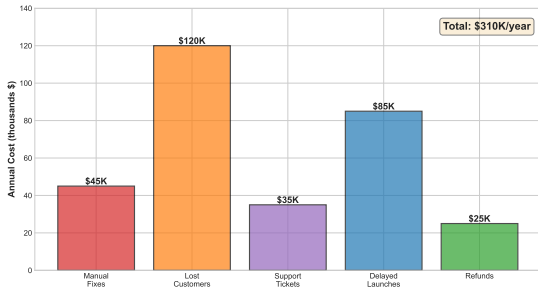- Production deployment

**Parts 4-5: Design & Practice**
- UX for reliable AI
- Workshop exercise
- Best practices
- Key takeaways

Making AI systems production-ready and trustworthy

Cost of Unreliable AI Outputs
Annual Impact Per 1000 Users

## $310K Per Year

### Impact Areas:

- Manual error correction
- Customer churn from mistakes
- Support ticket overload
- Delayed product launches
- Refunds and compensation

Per 1000 users - typical AI-powered service

## The 80% Problem: Why Most AI Projects Fail

### The Gap

- 80% of AI projects never reach production
- Prototypes work in demos, fail in reality
- Unpredictable outputs
- No integration path
- Cannot handle errors

Week 6: Generated creative content
Week 8: Make it reliable and usable

### The Solution

- Structured outputs
- JSON schema validation
- Error handling
- Production architecture
- Monitoring and testing

**Result:** Prototype → Production-ready MVP

Moving from creative exploration to reliable deployment

# When AI Goes Wrong: Real Examples

## E-commerce Chatbot

- Generated wrong pricing
- Promised impossible discounts
- Gave conflicting product info

Impact:
$45K in honored mistakes
2,300 confused customers
Brand damage

## Form Filling AI

- Inconsistent field extraction
- Mixed up phone/email
- Lost required data

Impact:
40% forms required manual fix
3 hours/day staff time
Customer frustration

## Report Generator

- Formatting varied wildly
- Missing key sections
- Unstructured data

Impact:
Reports unusable as-is
Lost automation benefits
Manual reconstruction

Common failure pattern: Unstructured outputs in structured contexts

# Structured vs Unstructured Outputs

## Unstructured Output

The restaurant was amazing! I'd give it 5 stars. Great food quality and service was excellent. Price was moderate around $30 per person.

## Structured Output (JSON)

```
{
  "rating": 5,
  "food_quality": 5,
  "service": 5,
  "price_level": "moderate",
  "avg_price_per_person": 30,
  "recommended_for": ["date", "friends"]
}
```

**Problems:**

- No standard format

- Requires parsing

- Error-prone extraction

- No validation

**Benefits:**

- Standard JSON format

- Direct integration

- Type validation

- Reliable parsing

Structured outputs enable reliable automation and integration

## When Do You Need Structured Outputs?

**Use Structured Outputs:**
- Database integration
- API responses
- Form filling
- Data extraction
- Automated workflows
- Multi-step processing
- Validation requirements
- Consistent formatting

When reliability matters more than creativity

Most production AI systems need structure for reliability

**Use Unstructured Text:**
- Creative writing
- Content generation
- Explanations
- Brainstorming
- Conversational responses
- Marketing copy
- Storytelling

When creativity matters more than structure

# What Makes AI Production-Ready?

**Technical Requirements**
- Consistent output format
- Schema validation
- Error handling
- Retry logic
- Monitoring
- Logging
- Performance SLAs
- Cost optimization

**Business Requirements**
- 95%+ success rate
- < 2 second response time
- Graceful degradation
- User trust
- Compliance
- Audit trails
- ROI positive
- Scalable

Structured outputs are the foundation for meeting these requirements

Production readiness requires reliability, not just functionality

# What You'll Master This Week

**Technical Skills**
1. Design JSON schemas for AI outputs
2. Implement function calling (OpenAI/Anthropic)
3. Write prompts for structured generation
4. Build validation pipelines
5. Handle errors gracefully
6. Deploy to production
7. Monitor system health

**Design Skills**
1. Create UX for AI features
2. Build trust through consistency
3. Design error recovery flows
4. Human-in-the-loop patterns
5. Progressive enhancement
6. Accessibility considerations

By the end: Transform prototypes into production MVPs

Practical skills for building real AI products

# Innovation Impact: Speed to Market

## Without Structured Outputs

- Prototype looks good
- Integration takes weeks
- Constant manual fixes
- Cannot scale
- User complaints
- Team loses confidence

Timeline: 6-12 weeks prototype → production
Success rate: 20%

## With Structured Outputs

- Prototype integrates directly
- Validation catches errors
- Automated workflows
- Scales to thousands
- Reliable user experience
- Team ships with confidence

Timeline: 1-2 weeks prototype → production
Success rate: 85%

Structured outputs dramatically reduce time-to-market

# Evolution of AI Reliability

**2020-2022**
Text Generation Era
- GPT-3 creative outputs
- Unstructured text
- Manual parsing required
- Low reliability
- Demo-only quality

**2023**
Function Calling Era
- OpenAI function calling
- JSON mode
- Structured outputs
- 90%+ reliability
- Production-ready

**2024-2025**
Reliable AI Era
- Native structured output
- Schema enforcement
- 99% reliability
- Enterprise-grade
- Mainstream adoption

We're at the tipping point: AI becomes truly reliable

The transition from creative tools to production systems

# Foundation Summary: Key Principles

## Core Concepts

1. Reliability is the production bottleneck
2. Structured outputs solve 80% gap
3. JSON schemas define contracts
4. Validation catches errors early
5. Monitoring ensures quality

**Remember:**
Creativity for exploration
Structure for production

## Success Metrics

- 95%+ success rate
- < 2s response time
- Zero manual parsing
- Direct integration
- User trust

**Next Steps:**
Learn the techniques to achieve this reliability

Part 2: Techniques for making AI outputs reliable

## JSON Schema Example

*Restaurant Review Validation*

Type constraints

Value validation

```json
{
  "type": "object",
  "properties": {
    "rating": {
      "type": "integer",
      "minimum": 1,
      "maximum": 5
    },
    "food_quality": {
      "type": "integer",
      "minimum": 1,
      "maximum": 5
    },
    "price_level": {
      "type": "string",
      "enum": ["cheap", "moderate", "expensive"]
    },
    "recommended_for": {
      "type": "array",
```

## Prompt Engineering Patterns: Success Rate & Consistency



Legend: Success Rate (%), Consistency (%)

| Pattern | Success Rate (%) | Consistency (%) |
|---|---|---|
| Basic Prompt | 72% | 68% |
| Role-Based | 81% | 79% |
| Step-by-Step | 88% | 91% |
| Few-Shot | 92% | 93% |
| Chain-of-Thought | 95% | 96% |

More structured prompts yield more consistent outputs

# Five Prompt Patterns Explained

**1. Basic Prompt**
"Extract data from this review"
Success: 72%

**2. Role-Based**
"You are a data extraction expert.  Extract..."
Success: 81%

**3. Step-by-Step**
"1.  Read review 2.  Identify rating 3.  Extract..."
Success: 88%

**4. Few-Shot**
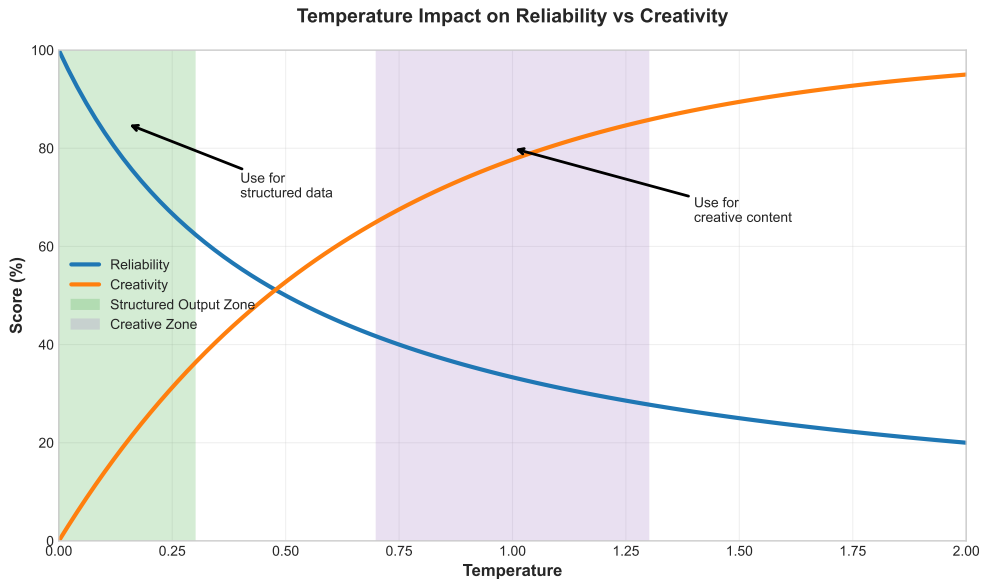Provide 2-3 examples
Success: 92%

**5. Chain-of-Thought**
"Think through each field.  Explain your reasoning..."
Success: 95%

Combine patterns for best results: Role + Few-Shot + CoT = 97%

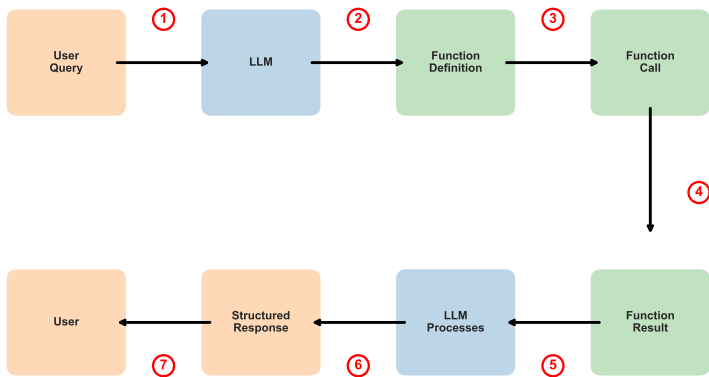Pattern selection depends on complexity and requirements

Temperature Impact on Reliability vs Creativity

**Function Calling Flow Architecture**

# Function Calling vs Tool Use: What's the Difference?

## Function Calling
OpenAI, Google

- Model generates function call
- You execute the function
- Return results to model
- Model processes response

**Best for:**

- Structured data extraction
- API integrations
- Multi-step workflows

## Tool Use
Anthropic Claude

- Model requests tool
- Same pattern, different API
- More explicit tool definitions
- Designed for agents

**Best for:**

- Agent systems
- Complex tool chains
- Interactive workflows

Both achieve structured outputs - choose based on your LLM provider

Conceptually similar, API differences only

# Chain-of-Thought: Improving Reasoning

**Without CoT**
`Extract: {rating: 3, price: "moderate"}`

Problems:

- No reasoning visible
- Hard to debug errors
- Inconsistent logic
- Cannot verify

**With CoT**
`Reasoning: "Customer mentions 'okay food' suggesting 3/5 stars. They say '$25 per person' which is moderate range."`
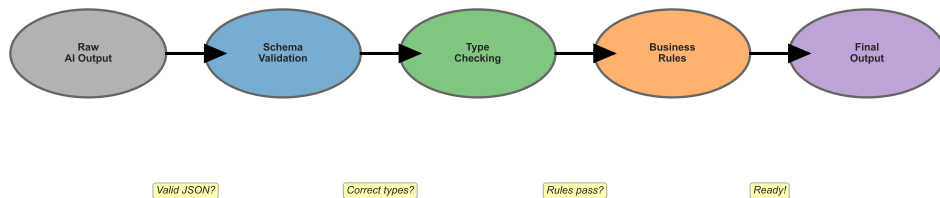`Extract: {rating: 3, price: "moderate"}`

Benefits:

- Reasoning traceable
- Easier debugging
- More consistent
- Verifiable logic

CoT improves accuracy by 5-15% for complex extractions

**Multi-Stage Validation Pipeline**



Layer validations to catch different types of errors

# Three Layers of Validation

### 1. Schema Validation

- Valid JSON?
- All fields present?
- Correct types?
- Within ranges?

Tools:
JSON Schema
Pydantic
TypeScript types

### 2. Business Rules

- Logical consistency?
- Cross-field validation?
- Domain constraints?
- Edge cases?

Example:
If rating = 5
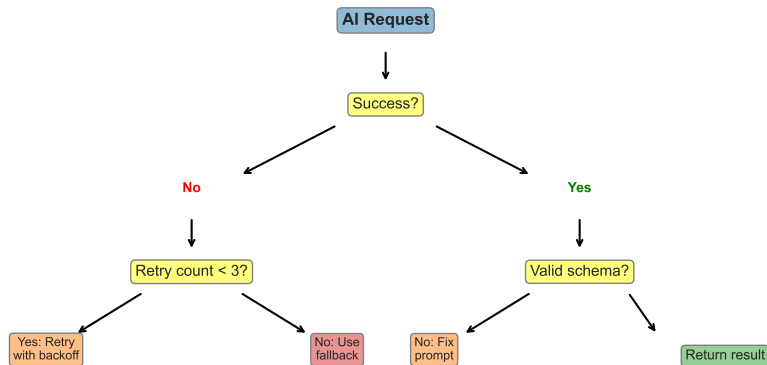then sentiment cannot be negative

### 3. Confidence Checks

- Model confidence score?
- Ambiguous input?
- Unusual values?
- Human review needed?

Action:
< 70% confidence
→ Flag for review

Each layer catches different failure modes

**Error Handling Decision Tree**



**Fallback Options:**

## Technique Selection Guide

| Technique | Reliability | Speed | Best For |
|---|---|---|---|
| Basic Prompt | 70-80% | Fast | Simple extraction |
| Role + Steps | 85-90% | Fast | Medium complexity |
| Few-Shot | 90-95% | Medium | Consistent format |
| Chain-of-Thought | 95-97% | Slow | Complex reasoning |
| Function Calling | 95-99% | Fast | Structured APIs |
| Multi-Validation | 98-99% | Medium | Critical data |

**Recommended:** Function calling + Few-shot + Validation
**Result:** 98%+ reliability at reasonable speed

Combine techniques for production-grade reliability

## OpenAI Function Calling: Code Example

```
functions = [{
      "name":  "extract_review",
      "description":  "Extract data",
      "parameters":  {
            "type":  "object",
            "properties":  {
                    "rating":  {"type":  "integer"},
                    "price":  {"type":  "string"}
            }
      }
}]

response = openai.ChatCompletion.create(
      model="gpt-4",
      messages=[...],
      functions=functions
)
```

OpenAI handles JSON schema validation internally

**Key Points:**

- Define schema upfront
- Model decides to call function
- Returns structured JSON
- Validates automatically

**Benefits:**

- Native validation
- Type-safe
- No parsing needed
- 95%+ reliability

## Anthropic Tool Use: Alternative Approach

```
tools = [{
      "name":  "extract_review_data",
      "description":  "Extract structured data",
      "input_schema":  {
              "type":  "object",
              "properties":  {
                      "rating":  {...
              },
              "required": ["rating"]
      }
}]

message = anthropic.messages.create(
      model="claude-3-opus",
      tools=tools,
      messages=[...]
)
```

**Differences:**

- input_schema vs parameters
- More explicit tool definitions
- Designed for multi-tool agents

**Same Result:**

- Structured JSON output
- Type validation
- High reliability

Choose based on your LLM provider - both work well

# Pydantic: Type-Safe Python Validation

```python
from pydantic import BaseModel
class Review(BaseModel):
        rating: int
        food_quality: int
        price_level: str

        @validator('rating')
        def check_rating(cls, v):
                if v < 1 or v > 5:
                        raise ValueError("1-5 only")
                return v

review = Review(**ai_output)
```
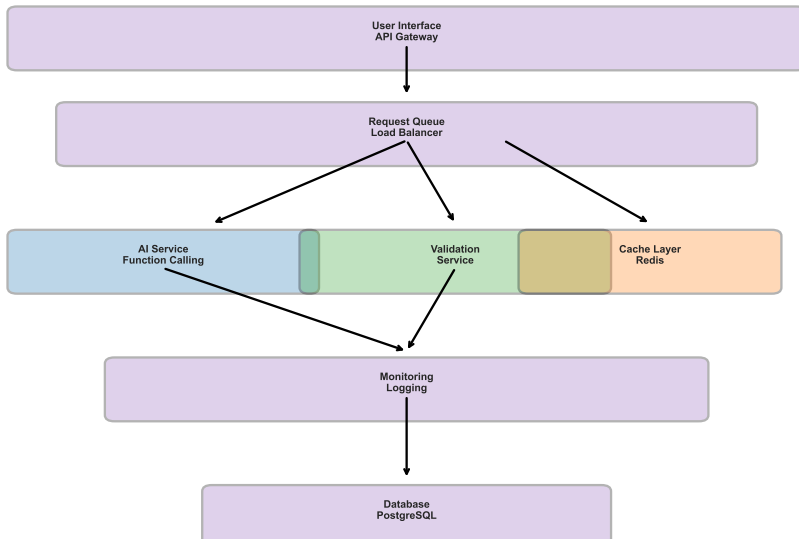
**Benefits:**
- Automatic type checking
- Custom validators
- Clear error messages
- IDE autocompletion
- JSON schema generation

**Production Ready:**
- Catches errors immediately
- Prevents bad data
- Self-documenting code

Pydantic is the standard for Python API validation

**Production Architecture for Structured AI**



User Interface
API Gateway

Request Queue
Load Balancer

AI Service
Function Calling

Validation
Service

Cache Layer
Redis

Monitoring
Logging

Database
PostgreSQL

## Graceful Error Handling Pattern

```
def extract_with_fallback(text, retries=3):
      for attempt in range(retries):
              try:
                      result = ai_extract(text)
                      if validate(result):
                              return result
                      else:
                              log_validation_failure(result)
              except APIError:
                      if attempt < retries - 1:
                              time.sleep(2 ** attempt) # Exponential backoff
                              continue

      # All retries failed - use fallback
      return rule_based_fallback(text)
```
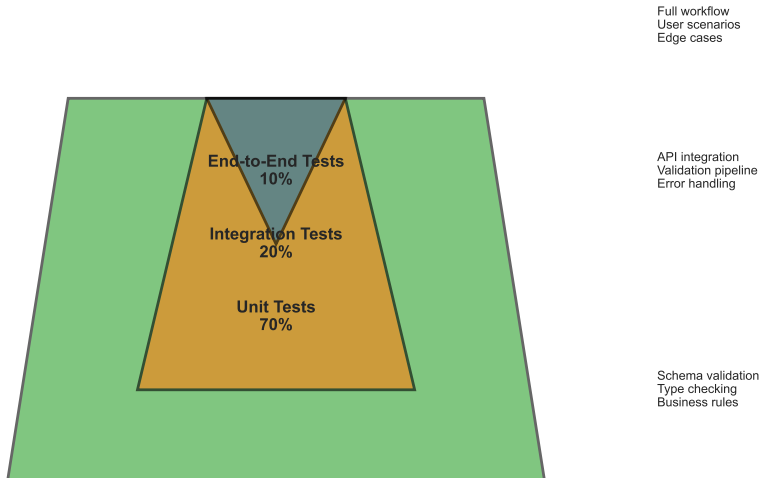
**Key Components:**

- Retry with exponential backoff
- Validation checks
- Logging for debugging
- Rule-based fallback
- Never return invalid data

Production systems need multiple fallback layers

## Testing Pyramid for Structured AI
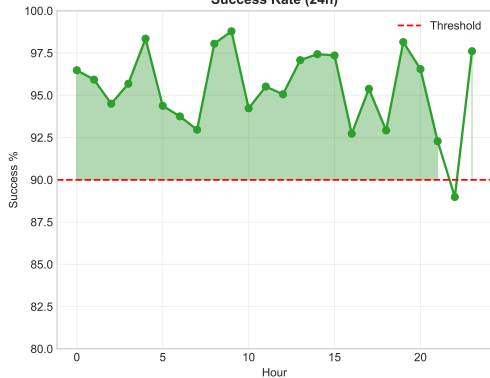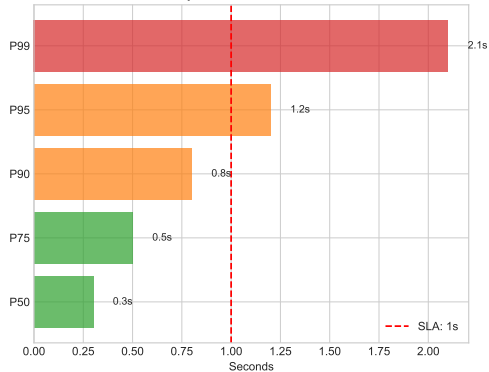


End-to-End Tests
10%

Integration Tests
20%

Unit Tests
70%

Full workflow
User scenarios
Edge cases

API integration
Validation pipeline
Error handling

Schema validation
Type checking
Business rules

Production Monitoring Dashboard

# Production Deployment Checklist

**Before Deployment**
- ☐ Schema defined and documented
- ☐ Validation tests pass 100%
- ☐ Error handling implemented
- ☐ Retry logic tested
- ☐ Fallback system works
- ☐ Logging configured
- ☐ Monitoring dashboards ready
- ☐ Alerts configured
- ☐ Load tested at 10x volume

**After Deployment**
- ☐ Success rate > 95%
- ☐ P95 latency < 2s
- ☐ Error rate < 2%
- ☐ Cost within budget
- ☐ No manual interventions needed
- ☐ User feedback positive
- ☐ Documentation updated
- ☐ Team trained
- ☐ Runbook created

Don't skip the checklist - it prevents production fires

Production readiness requires careful verification

# Optimization Strategies

### Token Reduction
- Shorter prompts
- Remove examples after tuning
- Compress context
- Use smaller models when possible

Impact:
50% cost reduction
30% faster

### Caching
- Cache identical requests
- 1-hour TTL
- Redis for speed
- Cache hit rate > 40%

Impact:
70% cost reduction
10x faster

### Batching
- Process multiple items together
- Async processing
- Queue management
- Batch size 10-50

Impact:
40% cost reduction
Better throughput

Optimization can reduce costs by 60-80% while maintaining quality

# Implementation Summary: Key Takeaways

## Core Implementation

1. Use function calling or tool use
2. Validate with Pydantic or similar
3. Implement retry + fallback
4. Add comprehensive logging
5. Monitor everything

### Production Requirements:

- 95%+ success rate
- < 2s P95 latency
- Graceful degradation
- Cost optimized

## Common Mistakes to Avoid

- No validation layer
- Single point of failure
- No error logging
- No monitoring
- Skipping testing
- No fallback plan
- Ignoring costs

**Next:** Design UX patterns for reliability

Part 4: Creating user experiences with structured AI

# UX Patterns for Reliable AI

## Progressive Loading

| Request sent |
|---|

| Processing... |
|---|

| Validating... |
|---|

| Ready! |
|---|

## Show Confidence

| High: 95% |
|---|

| Medium: 75% |
|---|

| Low: 45% |
|---|

*Please review*

## Graceful Error Recovery

Error Occurred

## Human-in-the-Loop

AI Suggestion

# Progressive Enhancement: Start Simple, Add AI

## The Pattern
1. Start with manual form
2. Add AI suggestions
3. User reviews and edits
4. Final submit

### Why It Works:
- User stays in control
- AI failures don't block
- Trust builds gradually
- Works without AI

Never make AI a single point of failure

## Example: Form Filling
1. User uploads invoice
2. AI extracts fields
3. Shows in editable form
4. User corrects mistakes
5. Saves valid data

### Result:
- 90% time saved
- 100% accuracy
- User confident

# Communicating AI Processing

## Stage-by-Stage Feedback

1. "Analyzing document..."
2. "Extracting data..."
3. "Validating fields..."
4. "Ready for review!"

### User Benefits:

- Knows what's happening
- Expected wait time
- Can cancel if needed
- Reduces anxiety

## Progress Indicators

- Spinner for $< 2s$
- Progress bar for 2-10s
- Stage labels for $> 10s$
- Time estimates when available

### What NOT to Do:

- Blank screen
- Generic "Loading..."
- No cancel option
- False progress bars

Clear feedback builds trust during AI processing

# User-Friendly Error Messages

**Bad Error Messages**
```
Error:  Schema validation failed at line 42
API returned 500
Unexpected token in JSON
```

Problems:

- Technical jargon
- No action suggested
- Scary and confusing
- User feels helpless

**Good Error Messages**
```
We couldn't process this document.  Please try:
```

- Upload a clearer image
- Enter data manually instead
- Contact support if this persists

Features:

- Plain language
- Actionable steps
- Alternative paths
- Reassuring tone

Error messages should help, not frustrate

## When to Show Confidence

- High-stakes decisions
- Ambiguous inputs
- User needs assurance
- Learning/training scenarios

## How to Display:

- Color coding (green/yellow/red)
- Percentage ("85% confident")
- Stars or bars
- Textual ("High confidence")

Confidence scores help users make informed decisions

## Confidence-Based Actions

| Confidence | Action |
| --- | --- |
| > 95% | Auto-accept |
| 80-95% | Suggest, allow edit |
| 60-80% | Show for review |
| < 60% | Request manual entry |

## Benefits:

- Appropriate review level
- User knows when to check
- Builds calibrated trust

# Human-in-the-Loop Patterns

## Three Levels of Human Control

1. High automation — AI decides, human monitors
2. Shared control — AI suggests, human approves
3. Human primary — AI assists, human decides

Choose based on:
- Risk level
- AI confidence
- User expertise
- Task complexity

Give users control appropriate to the task risk

## Example: Data Review Interface

**AI Extraction:**
- Shows extracted data
- Highlights low confidence
- Inline editing
- Accept/reject/edit options

**User Actions:**
- Quick accept if all good
- Edit specific fields
- Reject and re-extract
- Manual entry if AI fails

# Structured Input/Output UX

## Smart Form Filling

**User uploads document**

↓

**AI extracts fields**

↓

**Shows in form with indicators:**

- Green check: High confidence
- Yellow warning: Please review
- Red X: Couldn't extract

↓

**User edits as needed**

↓

**Validates before submit**

## Key UX Features

- Pre-filled, not read-only
- Clear confidence indicators
- Easy inline editing
- Field-level validation
- Show original source
- Undo/redo
- Save draft
- Skip AI option

**Result:**
90% time saved
User stays in control

Make AI suggestions obvious but easy to override

## Accessibility with Structured AI

### Why Structured Outputs Help
- Predictable format
- Screen reader friendly
- Keyboard navigation
- Clear structure
- Consistent patterns
- Alt text generation
- Semantic HTML

**Benefits:**
- WCAG 2.1 compliance easier
- Better for all users
- Legal requirements met

Structured data makes accessible AI easier to build

### Implementation Tips
- Use semantic elements
- ARIA labels for AI status
- Announce confidence levels
- Keyboard shortcuts
- Skip to error
- Focus management
- High contrast mode
- Text alternatives

**Example:**
```
<div role="status" aria-live="polite">
    AI extracted 8 of 10 fields
</div>
```

# Design for Trust

## Trust Through Consistency

- Predictable behavior
- Clear capabilities
- Honest about limits
- Graceful failures
- User stays in control

### Trust Builders:

- Show confidence scores
- Explain AI decisions
- Easy to override
- Consistent patterns
- No surprises

## Trust Destroyers:

- Inconsistent outputs
- Hidden AI decisions
- No way to correct
- Mysterious errors
- Overconfident claims
- Blocking failures
- No human override

### Golden Rule:
Underpromise and overdeliver

Trust is earned through consistent, reliable behavior

# Design Framework: Key Principles

## Core Principles

1. User always in control
2. Progressive enhancement
3. Clear feedback
4. Graceful degradation
5. Accessibility first
6. Build trust through consistency

## Structured AI Advantages:

- Predictable UI
- Easier to verify
- Clear error states
- Consistent patterns

## Checklist

- ☐ Works without AI
- ☐ Shows confidence
- ☐ Easy to edit
- ☐ Clear error messages
- ☐ Loading states
- ☐ Keyboard accessible
- ☐ Screen reader tested
- ☐ No blocking failures
- ☐ Cancel option
- ☐ User can override

Next: Put it all into practice with a workshop

Part 5: Hands-on workshop and best practices

# Workshop: Restaurant Review Intelligence System

**Your Challenge**
Build a system that extracts structured data from unstructured restaurant reviews.

**Why This Matters:**
- Real-world problem
- Applies all Week 8 concepts
- Production-ready skill
- Portfolio project

**Success Criteria:**
- 90%+ extraction accuracy
- Valid JSON output
- Handles errors gracefully

Complete, working system that extracts structured data reliably

**What You'll Build**
1. JSON schema definition
2. Extraction prompt
3. Function calling implementation
4. Validation pipeline
5. Error handling
6. Testing suite

Time: 60 minutes
Deliverable: Python notebook
Dataset: 1,000 reviews provided

## Workshop Dataset: 1,000 Restaurant Reviews

### Data Format
```
review_id:  1234
text:  "Amazing food!  The service was excellent..."
verified:  true
```

### Characteristics:
- 100-500 words per review
- Mix of positive/negative
- Various writing styles
- Different detail levels
- Some ambiguous cases

### Extract These Fields:
**Required:**
- overall_rating (1-5)
- food_quality (1-5)
- service_quality (1-5)
- price_level (cheap/moderate/expensive)

**Optional:**
- ambiance_rating (1-5)
- top_3_themes (array)
- recommended_for (array)

Dataset includes 100 human-labeled examples for validation

# Step-by-Step Implementation Guide

## Phase 1: Schema (15 min)
1. Define JSON schema
2. Add type constraints
3. Set value ranges
4. Mark required fields
5. Test with sample data

## Phase 2: Prompt (15 min)
1. Write extraction prompt
2. Add role definition
3. Include examples
4. Test on 5 reviews
5. Iterate to improve

## Phase 3: Implementation (20 min)
1. Set up function calling
2. Add validation layer
3. Implement error handling
4. Test on 50 reviews
5. Fix common failures

## Phase 4: Validation (10 min)
1. Run on 100 labeled examples
2. Calculate accuracy
3. Analyze failure cases
4. Document results

Starter notebook provided with code templates

# Testing & Validation Approach

## Unit Tests
- Schema validation works?
- Type checking catches errors?
- Business rules enforced?
- Edge cases handled?

## Integration Tests
- Full pipeline works?
- Error handling triggers?
- Retry logic functions?
- Fallback activates?

## Accuracy Metrics
- Field-level accuracy
- Overall match rate
- Confidence calibration
- Error type distribution

## Success Thresholds:
- Rating extraction: 95%+
- Price level: 90%+
- Themes: 85%+
- Overall system: 90%+

Compare your results against human-labeled ground truth

# Analyzing Your Results

## What to Analyze

1. Accuracy by field
2. Common error patterns
3. Confidence vs accuracy
4. Processing time
5. Cost per review
6. Edge case handling

### Questions to Ask:

- Which fields fail most?
- Why did specific cases fail?
- Is confidence score reliable?
- What patterns emerge?

Iteration is key - expect 2-3 refinement cycles

## Iteration Strategies

**If accuracy $< 90\%$:**

- Add more examples to prompt
- Refine schema constraints
- Improve error handling
- Lower temperature
- Try chain-of-thought

**If too slow:**

- Remove unnecessary steps
- Use smaller model
- Add caching
- Batch process

**Best Practices for Structured AI**

## Design

- Define clear JSON schema

- Document required vs optional fields

- Use enums for constrained values

- Include examples in schema

## Implementation

- Set temperature to 0-0.3

- Use function calling when available

- Implement multi-stage validation

- Add retry logic with backoff

## Testing

- Unit test schema validation

## Resources for Structured AI Development

### Libraries & Tools
**Python:**

- Pydantic - Validation
- OpenAI SDK - Function calling
- Anthropic SDK - Tool use
- JSON Schema - Definitions
- pytest - Testing

**Monitoring:**

- Datadog, New Relic
- Weights & Biases
- LangSmith

### Documentation

- OpenAI Function Calling Guide
- Anthropic Tool Use Tutorial
- Pydantic Documentation
- JSON Schema Validator
- Course handouts (3 levels)

### Practice Datasets:

- Restaurant reviews (today)
- Invoice extraction
- Customer support tickets
- Product descriptions

All resources linked in course materials

# Week 8 Key Takeaways

## Core Concepts

1. Structured outputs enable production AI
2. 80% of AI projects fail without reliability
3. JSON schemas define clear contracts
4. Validation catches errors early
5. Multiple techniques combine for 98%+ reliability

## Technical Skills:

- Function calling
- Pydantic validation
- Error handling
- Testing strategies
- Production deployment

## Design Skills:

- Progressive enhancement
- Confidence display
- Human-in-the-loop
- Error recovery UX
- Trust-building patterns

## Remember:

- Creativity for exploration
- Structure for production
- User always in control
- Trust through consistency

You can now build production-ready AI systems!

Next weeks: Testing, validation, and optimization