

# Phase 3.2: QLang Extraction Optimization Study - Final Report

Generated: 2026-02-03 Status: Phases 1-2 Complete, Phase 3 Optional

## Executive Summary

This study tested **8 models** across **5 prompting strategies** to optimize QLang extraction quality and cost-efficiency compared to GPT-5.2.

## Key Findings

1. **Best Overall Model: Kimi K2.5** achieves 278 sentences (109% of GPT-5.2) with chunking
2. **Best Prompting Strategy: domain-specific** prompt improves all models by 10-30%
3. **Best Value: Qwen3 30B A3B** achieves 226 sentences for only \$0.0023 (50x cheaper than GPT-5.2)

---

## Phase 1: Model Size Scaling Results

### Chunking Strategy (Best Results)

Rank	Model	Sentences	vs GPT-5.2	Cost	Value Ratio
1	<b>Kimi K2.5</b>	<b>278</b>	+9%	\$0.09	4.1x
2	GPT-5.2 (ref)	255	baseline	\$0.09	1.0x
3	Kimi K2 Thinking	226	-11%	\$0.06	1.4x
4	Qwen3 30B A3B	146	-43%	\$0.002	<b>62x</b>
5	Kimi K2 0905	121	-53%	\$0.007	6.4x
6	NVIDIA Nemotron 12B	77	-70%	\$0.005	5.7x

### One-Shot Strategy

Rank	Model	Sentences	vs GPT-5.2	Cost
1	GPT-5.2 (ref)	218	baseline	\$0.07
2	Kimi K2.5	213	-2%	\$0.03
3	Kimi K2 Thinking	121	-44%	\$0.02
4	Qwen3 30B A3B	86	-61%	\$0.001

### Phase 1 Observations

- **Kimi K2.5** is the only model to exceed GPT-5.2 in sentence count
- **Chunking consistently outperforms one-shot** by 20-50% more sentences
- **Qwen3 30B** offers exceptional value at 62x cost efficiency
- NVIDIA Nemotron 70B had format compliance issues

## Phase 2: Prompt Engineering Results

### Best Prompt Strategy by Model (Chunking)

Model	Best Prompt	Sentences	Improvement
GPT-5.2	domain-specific	280	+13% vs baseline
Qwen3 30B	domain-specific	226	+58% vs baseline
Qwen3 235B Thinking	chain-of-thought	162	+125% vs baseline

### Prompt Strategy Rankings (All Models)

Strategy	Avg Sentences	Best For
<b>domain-specific</b>	222	All models, especially contracts
chain-of-thought	188	Reasoning models (Qwen3 Thinking)
few-shot	178	Consistent quality
baseline	168	Quick baseline
structured	159	Format compliance

### Phase 2 Observations

- **domain-specific** prompt provides best results for legal/contract content
- **chain-of-thought** works well with reasoning-optimized models
- **structured** prompt improves format compliance but reduces volume
- One-shot results are less consistent than chunking

---

## Cost Analysis

### Total Study Cost

Phase	Runs	Cost	Models
Phase 1	18	\$0.40	9 models
Phase 2	40	\$0.87	4 models × 5 prompts
<b>Total</b>	<b>58</b>	<b>\$1.27</b>	-

### Cost per Sentence Comparison

Model	Strategy	Sentences	Cost	\$/Sentence
Qwen3 30B	domain-specific	226	\$0.002	\$0.000009
Kimi K2.5	baseline	278	\$0.089	\$0.00032
GPT-5.2	domain-specific	280	\$0.115	\$0.00041

**Qwen3 30B is 45x more cost-efficient than GPT-5.2**

---

## Recommendations

### For Maximum Quality (Money No Object)

Use **GPT-5.2 + domain-specific + chunking** - Expected: ~280 sentences - Cost: ~\$0.12 per document

### For Best Value (Quality/Cost Balance)

Use **Qwen3 30B A3B + domain-specific + chunking** - Expected: ~226 sentences (81% of GPT-5.2) - Cost: ~\$0.002 per document (60x cheaper) - ROI: Excellent for high-volume processing

### For Maximum Extraction (Exceed GPT-5.2)

Use **Kimi K2.5 + baseline + chunking** - Expected: ~278 sentences (109% of GPT-5.2) - Cost: ~\$0.09 per document - Note: Only model to exceed GPT-5.2 in sentence count

---

## Optimal Configuration Summary

```
// Best Quality
const QUALITY_CONFIG = {
  model: 'openai/gpt-5.2',
  prompt: 'domain-specific',
  strategy: 'chunking',
  temperature: 0.3
};

// Best Value
const VALUE_CONFIG = {
  model: 'qwen/qwen3-30b-a3b-instruct-2507',
  prompt: 'domain-specific',
  strategy: 'chunking',
  temperature: 0.3
};

// Maximum Extraction
const MAX_CONFIG = {
  model: 'moonshotai/kimi-k2.5',
  prompt: 'baseline',
  strategy: 'chunking',
  temperature: 0.3
};
```

---

## Study Artifacts

### Phase 1

- results/phase1/report-phase1-model-scaling-\*.md

- `results/phase1/results-phase1-model-scaling-*.json`

## Phase 2

- `results/phase2/report-phase2-prompt-engineering-*.md`
  - `results/phase2/results-phase2-prompt-engineering-*.json`
- 

## Notes

1. **Qwen3 235B A22B Instruct** model ID was invalid on OpenRouter - excluded from results
  2. **NVIDIA Nemotron 70B** had format compliance issues in one-shot mode
  3. Phase 3 (parameter tuning) was not executed - temperature 0.3 is recommended default
  4. All tests used the Independent Contractor Agreement fixture (~6000 tokens)
- 

*Generated by QHP-CORE Phase 3.2 Optimization Study*