

# Foundations of Language Modeling: Case Study

*Building a Domain-Specific Language Model for Intelligent Customer Support Auto-Completion*

---

## 1. Industry Context and Business Problem

---

### Company Profile: Meridian Financial Technologies

**Meridian Financial Technologies** is a mid-size fintech company that provides digital banking infrastructure to 45 regional credit unions across the United States. Their platform handles approximately 2.3 million customer support interactions per month across chat, email, and in-app messaging channels.

Meridian employs 340 customer support agents who respond to inquiries ranging from simple balance checks to complex loan modification requests. The agents use an internal support platform that includes a knowledge base, ticket management, and a response drafting interface.

### The Business Problem

Meridian's customer support operation faces three critical challenges:

#### 1. Response Time Degradation

Average first-response time has increased from 4.2 minutes to 7.8 minutes over the past 18 months as ticket volume grew 40%. Customer satisfaction scores (CSAT) have dropped from 4.3/5.0 to 3.7/5.0 in the same period. Internal analysis shows that 62% of the response time is spent on the agent *typing* the response, not understanding the problem.

#### 2. Terminology Inconsistency

Financial services require precise language. Agents frequently use inconsistent terminology when discussing regulated topics -- for example, mixing up "annual percentage rate" and "interest rate," or using informal language for compliance-sensitive communications. A recent audit found terminology errors in 18% of customer-facing messages, exposing Meridian to regulatory risk.

#### 3. Onboarding Bottleneck

New agents take an average of 6 weeks to reach proficiency, primarily because they must learn Meridian's specific product terminology, compliance language, and response patterns. During

this ramp-up period, new agents require 3.2x more supervisor review than experienced agents, creating a bottleneck that limits hiring capacity.

## The Proposed Solution

Meridian's engineering team proposes building a **domain-specific language model** that provides intelligent auto-completion for customer support agents. Unlike generic auto-complete (which suggests common English phrases), this model must:

- Predict the next word/phrase using Meridian's specific financial terminology
- Suggest compliance-appropriate language for regulated topics
- Learn from the patterns of experienced agents to accelerate onboarding
- Run with sub-100ms latency for real-time suggestions

## Constraints

Constraint	Specification
Budget	\ 85,000 for initial development, 12,000/month ongoing compute
Latency	Auto-complete suggestions must appear within 100ms
Data	14 months of historical support transcripts (~3.2M messages)
Privacy	Customer PII must be stripped before model training
Deployment	Must run on Meridian's existing GPU infrastructure (4x NVIDIA A10G)
Compliance	Suggestions must never include financial advice or guarantees
Adoption	Target 60% agent adoption within 3 months of launch

## 2. Technical Problem Formulation

### From First Principles: What Are We Actually Building?

At its core, Meridian needs a system that, given a partially typed response and the customer's query as context, predicts the most likely next tokens. This is exactly the **language modeling** problem.

Let  $c = (c_1, c_2, \dots, c_m)$  be the customer's message (context) and  $r = (r_1, r_2, \dots, r_{t-1})$  be the agent's partially typed response. We want to compute:

$$P(r_t \mid c_1, \dots, c_m, r_1, \dots, r_{t-1})$$

This is conditional next-token prediction -- the same objective used to train GPT, but conditioned on both the customer query and the partial response.

## Architecture Decision: Why Not Just Use GPT?

The team evaluates three approaches:

### Option A: N-gram Model with Domain Vocabulary

Build a bigram/trigram model trained on historical support transcripts.

- Pros: Fast inference (~1ms), simple to deploy, fully interpretable
- Cons: Cannot capture long-range dependencies (e.g., referring back to the customer's loan number), severe sparsity on rare financial terms, no understanding of semantic similarity between terms

### Option B: Fine-Tuned GPT-2 (124M parameters)

Fine-tune an open-source GPT-2 model on Meridian's support corpus.

- Pros: Strong language understanding, handles long-range context, captures semantic relationships
- Cons: 124M parameters requires ~500MB memory, inference latency ~50-80ms on A10G, risk of generating hallucinated financial information

### Option C: Custom Small Transformer (8M parameters)

Train a purpose-built Transformer language model from scratch on Meridian's domain data.

- Pros: Optimized for domain vocabulary, small enough for <20ms inference, full control over training data and behavior, no risk of leaking pre-training knowledge
- Cons: Requires more engineering effort, smaller model capacity

The team selects **Option C** with an **N-gram fallback** (Option A) for edge cases where the Transformer produces low-confidence predictions.

## Loss Function and Training Objective

The model is trained with standard **cross-entropy loss** on next-token prediction:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=1}^T \log P_{\theta}(r_t \mid c, r_1, \dots, r_{t-1})$$

where  $\theta$  represents the model parameters.

The input to the model is a concatenation of the customer message and the agent's partial response, separated by a special token:

$$\text{input} = [\text{CLS}] \ c_1 \ c_2 \ \dots \ c_m \ [\text{SEP}] \ r_1 \ r_2 \ \dots \ r_{t-1}$$

The model only computes loss on the agent response tokens (not the customer message), since we are training to predict what the agent will write, not what the customer said.

## Evaluation Metrics

Metric	Definition	Target
<b>Perplexity</b>	$\exp(\mathcal{L})$ on held-out agent responses	< 15
<b>Top-1 Accuracy</b>	% of time the top suggestion matches the next token	> 35%
<b>Top-5 Accuracy</b>	% of time the correct token is in the top 5 suggestions	> 65%
<b>Keystroke Savings</b>	% reduction in keystrokes when agents use suggestions	> 40%
<b>Terminology Compliance</b>	% of suggestions using approved terminology	> 95%
<b>Inference Latency</b>	P99 latency for generating top-5 suggestions	< 100ms

## Baseline: Bigram Model

Before building the Transformer, the team establishes a baseline using a bigram model trained on the same corpus:

Metric	Bigram Baseline
Perplexity	142.3
Top-1 Accuracy	18.2%
Top-5 Accuracy	31.7%
Inference Latency	1.2ms

The bigram model is fast but produces poor suggestions because financial support language is highly contextual -- the right next word depends heavily on what was discussed earlier in the conversation, not just the immediately preceding word.

## 3. Implementation Notebook Structure

### 3.1 Data Preparation and Exploration

Load and explore Meridian's support transcript dataset. Strip PII using regex patterns and named entity recognition. Analyze vocabulary distribution, message lengths, and common response patterns.

```
# Load the support transcript dataset
# Format: (customer_message, agent_response, metadata)

import pandas as pd
import re
from collections import Counter

# TODO: Load the dataset from the provided CSV
# dataset = pd.read_csv("meridian_support_transcripts.csv")

# TODO: Implement PII stripping
# def strip_pii(text):
```

```
# """Remove customer PII from text."""
# text = re.sub(r'\b\d{4}[-\s]?d{4}[-\s]?d{4}[-\s]?d{4}\b', '[CARD_NUMBER]', text)
# text = re.sub(r'\b\d{3}-\d{2}-\d{4}\b', '[SSN]', text)
# text = re.sub(r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b', '[EMAIL]', text)
# return text
```

## 3.2 Tokenizer Construction

Build a domain-specific BPE tokenizer trained on Meridian's vocabulary. Financial terms like "APR," "amortization schedule," and "overdraft protection" should be single tokens.

```
# TODO: Train a BPE tokenizer on the support corpus
# from tokenizers import Tokenizer, models, trainers, pre_tokenizers

# tokenizer = Tokenizer(models.BPE())
# tokenizer.pre_tokenizer = pre_tokenizers.ByteLevel()
# trainer = trainers.BpeTrainer(vocab_size=8000, special_tokens=["[PAD]", "[CLS]", "[SEP]", "[UNK]"])
# tokenizer.train_from_iterator(all_texts, trainer=trainer)
```

## 3.3 N-gram Baseline Implementation

Implement the bigram and trigram baseline models from Notebook 1 of the teaching series. Compute perplexity and top-k accuracy on the validation set.

```
# TODO: Implement bigram baseline on Meridian's corpus
# Build count tables, compute probabilities, evaluate perplexity
```

## 3.4 Transformer Model Architecture

Build the custom 8M-parameter Transformer language model. Key design decisions: 6 layers, 8 heads, d\_model=256, context length=512 tokens.

```
# TODO: Implement the Meridian Transformer LM
# class MeridianLM(nn.Module):
#     def __init__(self, vocab_size=8000, d_model=256, num_heads=8,
#                   num_layers=6, max_len=512, dropout=0.1):
#         super().__init__()
#         # Token + position embeddings
#         # Transformer blocks with causal attention
#         # Output projection
#         pass
```

## 3.5 Training Pipeline

Implement the training loop with learning rate warmup, gradient clipping, and mixed-precision training. Track perplexity on train and validation sets.

```
# TODO: Training loop with:
# - AdamW optimizer (lr=3e-4, weight_decay=0.01)
# - Linear warmup over 1000 steps
# - Cosine annealing learning rate schedule
# - Gradient clipping at norm 1.0
# - Mixed precision (torch.cuda.amp)
# - Validation every 500 steps
```

### 3.6 Inference and Auto-Complete Engine

Build the real-time suggestion engine that takes partial input and returns top-k completions with confidence scores.

```
# TODO: Implement the auto-complete engine
# class AutoCompleteEngine:
#     def __init__(self, model, tokenizer, max_suggestions=5):
#         pass
#
#     def suggest(self, customer_message, partial_response, top_k=5):
#         """Return top-k next-token suggestions with probabilities."""
#         pass
#
#     def suggest_phrase(self, customer_message, partial_response, max_tokens=8):
#         """Return multi-token phrase suggestions using beam search."""
#         pass
```

### 3.7 Compliance Filter

Implement a post-processing filter that blocks suggestions containing financial advice, guarantees, or non-compliant terminology.

```
# TODO: Implement compliance filtering
# blocked_patterns = [
#     r"guarantee",
#     r"we promise",
#     r"you will (definitely|certainly|surely)",
#     r"your (money|funds|investment) (is|are) (safe|secure|protected)",
# ]
#
# approved_terms = {
#     "interest rate": "annual percentage rate (APR)",
#     "fee": "service charge",
#     "penalty": "assessed charge",
# }
```

### 3.8 Evaluation and A/B Testing Framework

Evaluate the model against the baseline on all metrics. Set up an A/B testing framework for production deployment.

```
# TODO: Full evaluation pipeline
# - Perplexity on held-out test set
# - Top-1 and Top-5 accuracy
# - Keystroke savings simulation
# - Terminology compliance rate
# - Latency benchmarking (P50, P95, P99)
# - Comparison table: Bigram vs Transformer
```

### 3.9 Production Deployment

Package the model for production deployment with ONNX export, batched inference, and monitoring.

```
# TODO: Production deployment preparation
# - Export model to ONNX for optimized inference
# - Implement KV-cache for efficient autoregressive generation
# - Set up latency monitoring and model quality dashboards
# - Implement gradual rollout with fallback to bigram model
```

---

## 4. Production and System Design Extension

---

### System Architecture

The production auto-complete system consists of three layers:

**Layer 1: Input Processing Pipeline** - Customer message is tokenized and encoded - Agent's partial response is concatenated with a separator token - PII detection runs in parallel (any detected PII triggers the input to be blocked from model logging)

**Layer 2: Model Inference** - The Transformer LM runs on a dedicated A10G GPU - KV-cache stores attention key/value pairs to avoid recomputation as the agent types each character - Batch inference: multiple agents' requests are batched together for GPU efficiency - Fallback: if P99 latency exceeds 100ms, the system falls back to the trigram model

**Layer 3: Post-Processing and Ranking** - Top-k token predictions are decoded back to text - Compliance filter removes prohibited suggestions - Terminology normalizer replaces non-standard terms with approved alternatives - Multi-token suggestions are generated via beam search (beam width = 3) - Final suggestions are ranked by:  $P(\text{token}) * \text{compliance\_score} * \text{diversity\_bonus}$

### Scaling Considerations

**Data Pipeline:** - New support transcripts are collected daily (PII-stripped) - Model is retrained weekly on the latest 6 months of data - A/B testing framework compares new model versions against the production model - Canary deployment: new models serve 5% of traffic for 48 hours before full rollout

**Model Scaling:** - Current model: 8M parameters,  $d_{\text{model}}=256$ , 6 layers - If domain grows: scale to 25M parameters ( $d_{\text{model}}=384$ , 8 layers) - If multi-language support needed: train separate tokenizers per language, share Transformer backbone

**Monitoring:** - Real-time dashboards: suggestion acceptance rate, latency, compliance violations - Drift detection: monitor vocabulary distribution shift (new products, new regulations) - Agent feedback loop: agents can flag bad suggestions, creating a negative training signal

### Expected Impact

Metric	Before	Target (6 months)
Avg response time	7.8 min	4.5 min
CSAT score	3.7/5.0	4.2/5.0
Agent onboarding time	6 weeks	3.5 weeks

Metric	Before	Target (6 months)
Terminology error rate	18%	4%
Agent throughput	28 tickets/day	42 tickets/day

## Cost Analysis

Component	Monthly Cost
GPU compute (1x A10G for inference)	\\$3,200
GPU compute (training, weekly)	\\$1,800
Data pipeline and storage	\\$1,500
Monitoring and logging	\\$800
<b>Total</b>	<b>\\$7,300/month</b>

With 340 agents handling 42 tickets/day (up from 28), the effective capacity increase is equivalent to hiring ~170 additional agents at an average cost of \\$4,200/month each. The ROI is approximately **98x** on infrastructure investment alone.

## Lessons Learned: From Theory to Production

- N-grams are still valuable as fallbacks.** The trigram model runs at 1ms latency and provides reasonable suggestions when the Transformer is under load or when the input is too short for meaningful attention patterns.
- Domain tokenization matters enormously.** Generic BPE tokenizers split "amortization" into subwords, reducing prediction accuracy. A domain-trained tokenizer that keeps financial terms intact improved top-1 accuracy by 8 percentage points.
- The sparsity problem from Notebook 1 reappears in production.** Rare but important financial terms (e.g., "forbearance agreement") appear too infrequently for the model to learn good predictions. The solution: augment training data with curated examples from the compliance team.
- Attention patterns reveal what the model learned.** Visualizing attention weights showed that the model learned to attend to specific entities in the customer message (account type, product name) when generating the response -- exactly the behavior the team wanted.
- Temperature tuning is critical for auto-complete.** Too low (0.3): suggestions are repetitive and predictable. Too high (1.0): suggestions are creative but sometimes nonsensical. The sweet spot for support auto-complete was temperature=0.5 with top-k=10 filtering.