

Call Me Maybe?

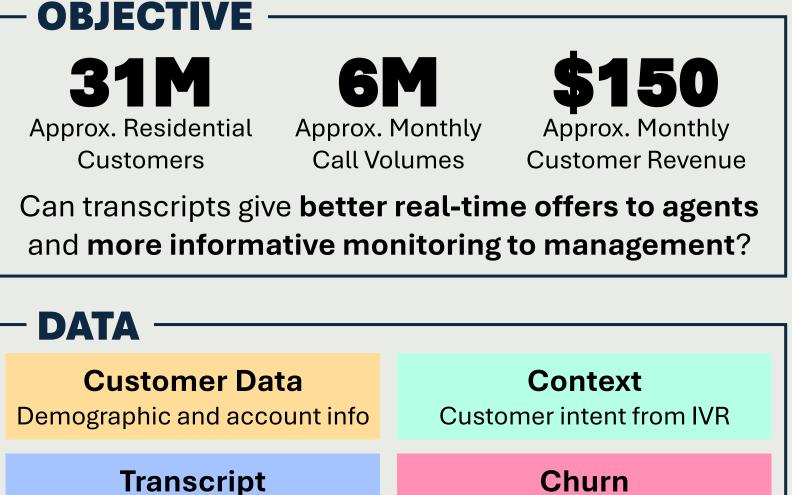
Predicting Churn in Real-time Using NLP & LLMs

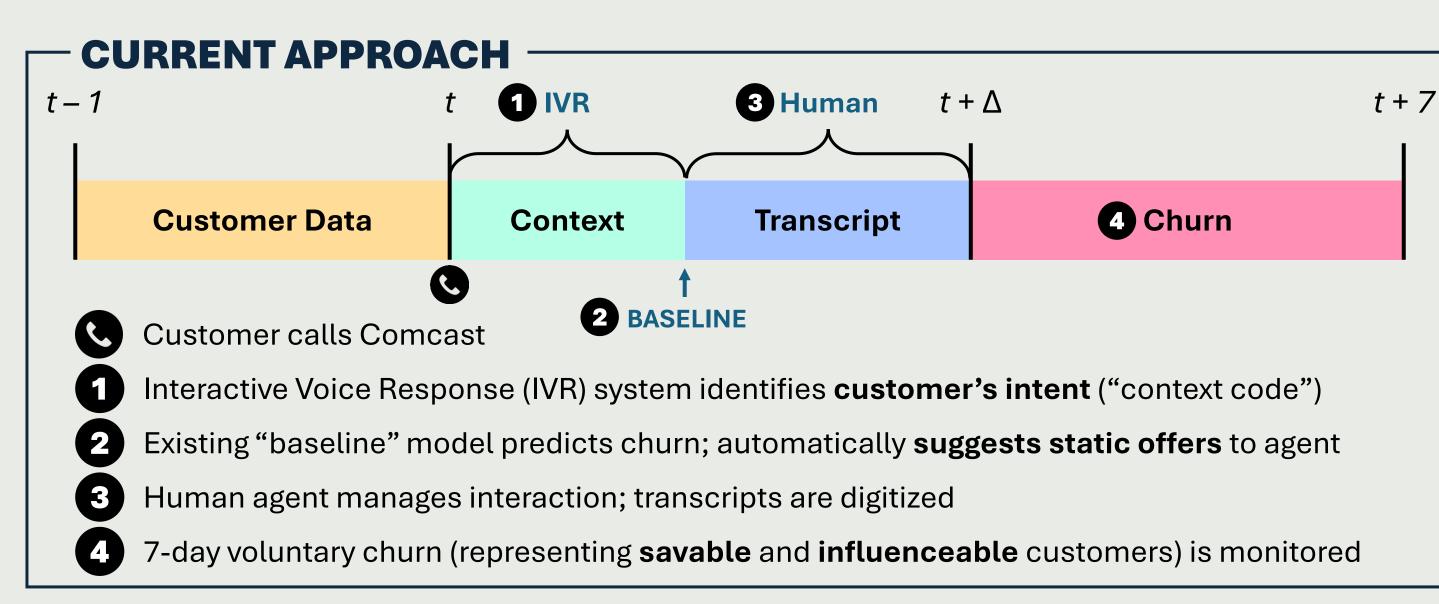
Comcast Project Sponsors: Benjamin Land, Ian Tongs, & Jacob Linder Faculty Advisor: James Butler

MBAn Students: Dilan SriDaran (dilan_s@mit.edu), Maxime Wolf (maximew@mit.edu)



PROBLEM STATEMENT





IMPACT

Digitized transcripts

REAL-TIME MODEL IMPLEMENTED AFTER 100 WORDS OF TEXT Potentially savable monthly revenue High-risk customers identified **Customer Data** Churn **Transcript** Context Can provide early, actionable, and dynamic insights and offers Latency and model **simplicity is a key constraint** for deployment Lower performance due to less information from transcripts

7-day voluntary churn flag

NON-REAL-TIME MODEL IMPLEMENTED AFTER CALL IS COMPLETE

High-risk customers identified

Potentially savable monthly revenue

Customer Data

Context

Transcript

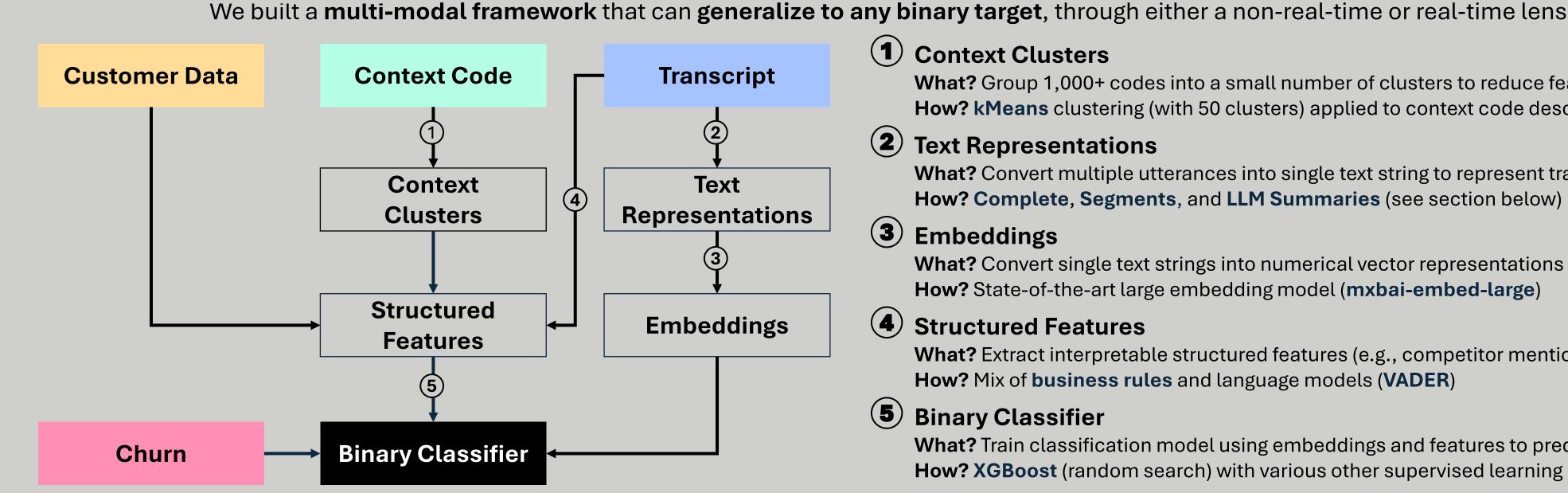
Churn

- Makes use of **all available information** from transcripts
- Latency and model complexity is not a constraint
- Less scope to act on predictions after a call

For internal security, results shown based on rounded assumptions and figures extracted from publicly released quarterly financials

APPROACH

GENERAL FRAMEWORK



Context Clusters

What? Group 1,000+ codes into a small number of clusters to reduce feature granularity How? kMeans clustering (with 50 clusters) applied to context code description embeddings

Text Representations

What? Convert multiple utterances into single text string to represent transcript How? Complete, Segments, and LLM Summaries (see section below)

Embeddings

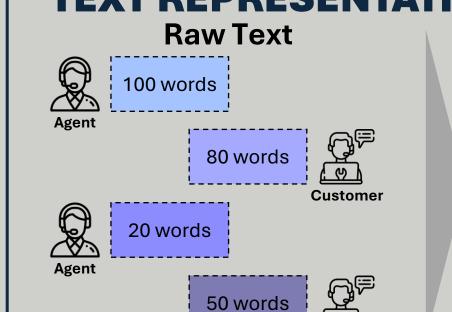
What? Convert single text strings into numerical vector representations **How?** State-of-the-art large embedding model (mxbai-embed-large)

Structured Features

What? Extract interpretable structured features (e.g., competitor mentions) from transcripts **How?** Mix of **business rules** and language models (**VADER**)

Binary Classifier

What? Train classification model using embeddings and features to predict 7-day churn How? XGBoost (random search) with various other supervised learning algorithms tested



(O)

Option 1: Complete Single text string.

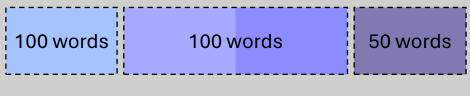
Single embedding.

250 words

- Computationally lightweight
- Embeddings often poor on long texts
- X Not applicable for real-time settings

Option 2: Segments

Multiple text strings. One embedding per string. Average of embeddings up to current segment.



- Computationally **lightweight**
- ✓ Suitable in both real- and non-real-time

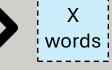
Option 3: LLM Summaries

Single text summary via LLM. Single embedding.



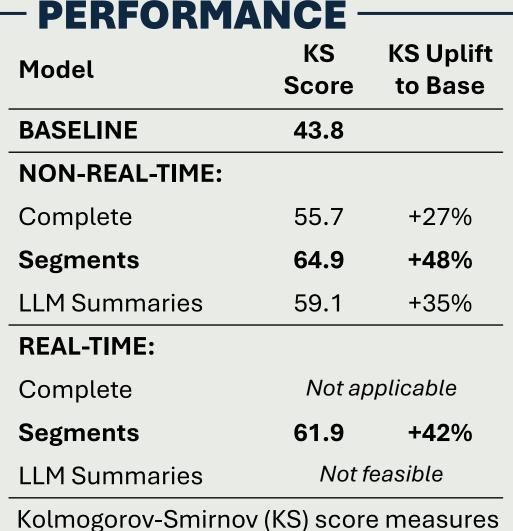




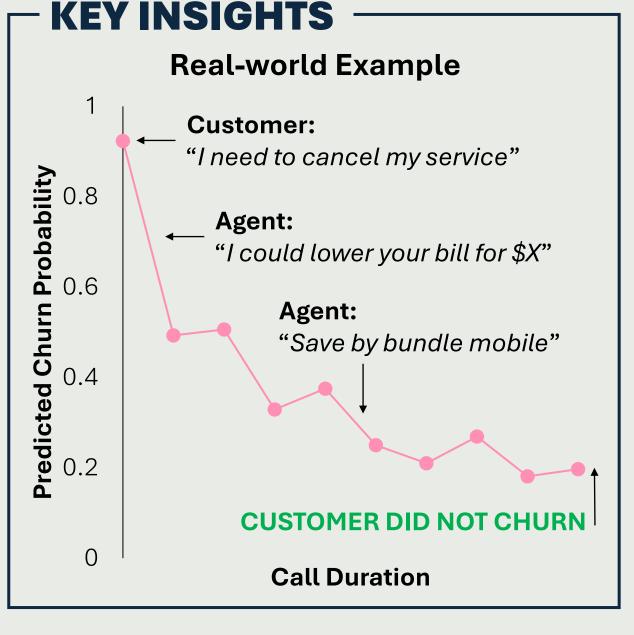


- ✓ Embeddings better on shorter texts
- Summaries computationally intensive
- ✓ Embeddings better on shorter texts
- ✓ Enhances model explainability

RESULTS



purity of separation between predicted classes; higher value is better



USE CASES



Real-time Offers

Re-run risk prediction early in call and revise offers; lower churn for high-risk customers and lower dilution for low-risk



Post-call Follow-ups Targeted follow-up calls to high-risk customers to ensure their issues were resolved and/or offer promotions

Transcript Mining

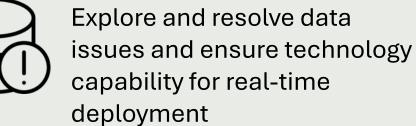


Mine transcripts using risk scores to understand churn drivers and competitive landscape

THE ROADMAP



Data Pipeline





Detailed Costing

Undertake costing including investments, LLM and compute costs, potential uplift, and dilution

Trial



Implement formal trial to explore effectiveness of programs and any unintended consequences