# BlendX : Complex Multi-intent Detection with Blended Patterns

Yejin Yoon, Jungyeon Lee, Kangsan Kim, Chanhee Park and Taeuk Kim

Yejin Yoon

Natural Language Processing Lab.,
Hanyang University.

HYU 한양대학교 HANYANG UNIVERSITY

Accepted by
# LREC-COLING 2024

# BlendX: Complex Multi-Intent Detection with Blended Patterns

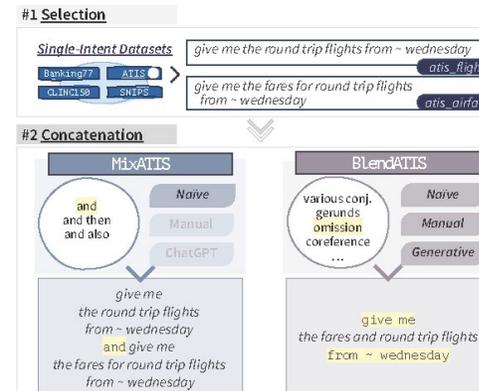## Anonymous submission

### Abstract

Task-Oriented Dialogue (TOD) systems typically suppose that a user utterance corresponds to a single intent. This assumption may be misaligned with real-world scenarios where users often express multiple intents simultaneously. While there is an emerging interest in Multi-Intent Detection (MID), existing in-domain datasets such as MixATIS and MixSNIPS have limitations in their formulation. To address these issues, we present BlendX, a suite of refined datasets featuring more diverse patterns than their predecessors, elevating both its complexity and difficulty. For dataset construction, we utilize both rule-based heuristics as well as a generative tool–OpenAI's ChatGPT—which is augmented with a similarity-driven strategy for utterance selection. To ensure the quality of the proposed datasets, we also introduce three novel metrics that assess statistical properties of an utterance related to word count, conjunction use, and pronoun usage. Extensive experiments on BlendX reveal that state-of-the-art MID models struggle with the challenges posed by the new datasets, highlighting the need to reexamine the current state of the MID field.

## 1. Introduction

The successful implementation of Task-Oriented Dialogue (TOD) systems begins with the precise recognition of user intents. By accurately discerning the queries embedded in user inputs and routing them to the relevant components, the systems can adeptly respond, thereby effectively fulfilling user requests. In general, such systems are constructed on the assumption that each user utterance correlates exclusively with a single intent, which often diverges from practical scenarios.

Contrary to the conventional setting, the task of **Multi-Intent Detection (MID)** presents a more nuanced and comprehensive challenge for TOD

# Contents

BlendX : Complex Multi-intent Detection with Blended Patterns

HYU 한양대학교
HANYANG UNIVERSITY

# PRE-REQUISITE

# Task-oriented Dialogue System

# Multi-intent Detection

HYU 한양대학교
HANYANG UNIVERSITY

# Task-oriented Dialogue System

play harry styles falling
and add it to favorite playlist

**Sound Signal**

**Text Input**

Speech Recognition

**NLU**
- Domain Identification
- User Intent Detection ✓
- Slot Filling

**Semantic Frame** ✓
PlayMusic, AddToPlaylist. ✓
song_title=falling, singer=harry styles, playlist_name=favorite

Now playing 'Falling' by Harry Styles.

**Text response**

**Response**
- Response Selection
- Response Generation **(NLG)**

**System Action / Policy**

confirm action

**Dialogue Management (DM)**
- Dialogue State Tracking
- Dialogue Policy

{call play music API}, {call add to playlist API}

Backend Action / Knowledge Providers

**Action**

- Help users achieve their specific goals
- Focus on understanding users, tracking states, and generating next actions

Natural Language Processing Lab.,
Hanyang University.

HYU 한양대학교 HANYANG UNIVERSITY

5

# Multi-intent Detection



- Input: 1 utterance
- Output: 1 or more label(s)

$if$ **threshold** $= 0.5$,

Legend:
- intent 0
- PlayMusic
- AddToPlaylist
- intent 1
- intent 2
- intent 3

Sigmoid

[CLS]  $T_1$  $T_2$  ...  $T_N$

**Pretrained Language Model (PLM)**
( e.g. BERT, RoBERTa, ELECTRA, …)

[CLS]  $Tok_1$  $Tok_2$  ...  $Tok_N$

**Tokenizer**

Since it is conversational speech,

omissions, contractions, and ungrammatical structures may occur.

play harry styles falling and add it to favorite playlist

Identify and respond to multiple intentions or requests within a single user utterance

# The Role of MID in ToD

According to a 2019 paper published by AWS AI, Amazon,
**52%**
**more than half** of its internal data utterances had multiple intentions.

amazon
Internal **ToD** Datasets

52%

Rashmi Gangadharaiah and Balakrishnan Narayanaswamy. Joint Multiple Intent Detection and Slot Labeling for Goal-Oriented Dialog. NAACL. 2019

# Introduction

# Problem States

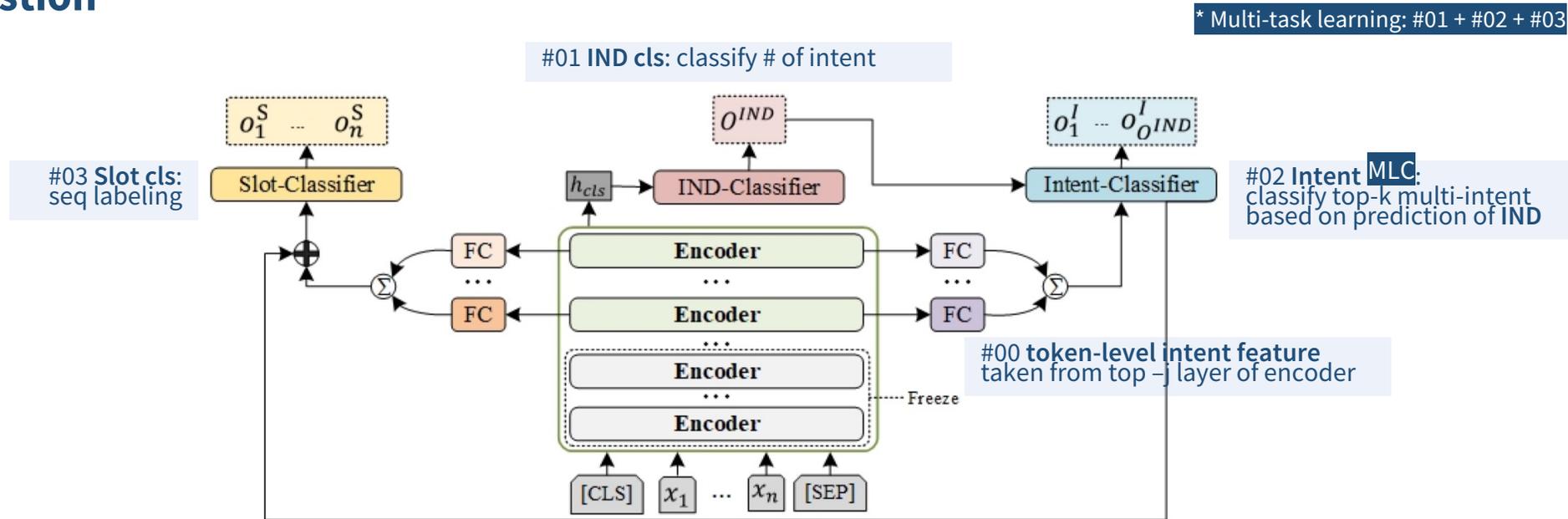# Background

# Related Works

- ## Literature Review
  - MID, Jointly-learning (w/ Slot Filling), MLC (Multi-Label Classification), Dataset

| Acronyms | Title | Authors | Released | Datasets | Categories | Date |
|---|---|---|---|---|---|---|
| CIBA | Multi-Point Semantic Representation for Intent Classification | Jinghan Zhang, et al. | AAAI2020 | CCL,CitySrv, ECOM,TELE | MID | 4/27 |
| AGIF | AGIF: An Adaptive Graph-Interactive Framework for Joint Multiple Intent Detection and Slot Filling | Libo Qin, et al. | EMNLP2020 Findings | MixATIS, MixSNIPS | Jointly-learning | 4/27 |
| GL-GIN | GL-GIN: Fast and Accurate Non-Autoregressive Model for Joint Multiple Intent Detection and Slot Filling | Fuxuan Wei, et al. | ACL-IJCNLP 2021 | MixATIS, MixSNIPS | Jointly-learning | 4/27 |
| MCT&ALR | Few-shot Learning for Multi-label Intent Detection | Yongkui Lai, et al. | AAAI2021 | StandfordLU, TourSG (DSTC-4) | MID w/ dynamic threshold | 4/27 |
| SDJN | Joint Multiple Intent Detection and Slot Filling via Self-distillation | Lisong Chen, et al. | ICASSP 2022 | MixATIS, MixSNIPS | Jointly-learning | 4/27 |
| ReLa-Net | Group is better than individual: Exploiting Label Topologies and Label Relations for Joint Multiple Intent Detection and Slot Filling | Bowen Xing, et al. | EMNLP2022 | MixATIS, MixSNIPS | Jointly-learning | 4/27 |
| AIK | Towards Multi-label Unknown Intent Detection | Yawen Ouyang,et al. | COLING2022 | MixSNIPS, MultiWOZ 2.3 | MID w/ out-of-scope | 4/27 |
| HBGL | Exploiting Global and Local Hierarchies for Hierarchical Text Classification | Ting Jiang, et al. | EMNLP2022 | WOS, NYT, RCV1-V2 | MLC w/ label semantics | 5/11 |
| Balanced LossNLP | Balancing Methods for Multi-label Text Classification with Long-Tailed Class Distribution | Yi Huang, et al. | EMNLP2021 | Reuters-21578, PubMed | MLC w/ class imbalance | 5/11 |
| MULTI-CONVFIT | Multi-Label Intent Detection via Contrastive Task Specialization of Sentence Encoders | Ivan Vulić, et al. | EMNLP2022 | MixATIS, NLU++ | MID w/ contrastive learning | 5/11 |
| KCOD | Watch the Neighbors: A Unified K-Nearest Neighbor Contrastive Learning Framework for OOD Intent Discovery | Yutao Mou, et al. | EMNLP2022 | Banking, CLINC, HWU64 | ID w/ contrastive learning | 5/11 |
| GISCo | Enhancing Joint Multiple Intent Detection and Slot Filling with Global Intent-Slot Co-occurrence | Mengxiao Song, et al. | EMNLP2022 | MixATIS, MixSNIPS | Jointly-learning | 5/11 |
| DialogUSR | DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection | Haoran Meng, et al. | EMNLP2022 Findings | *propose the datasets: DialogUSR | Dataset | 5/11 |

HYU 한양대학교 HANYANG UNIVERSITY

# Related Works ; Baseline

📄 **TFMN**) Chen et al. A Transformer-based Threshold-Free Framework for Multi-Intent NLU, COLING, 2022

- **Suggestion**



#01 **IND cls**: classify # of intent

\* Multi-task learning: #01 + #02 + #03

#03 **Slot cls**: seq labeling

#02 **Intent** MLC: classify top-k multi-intent based on prediction of **IND**

#00 **token-level intent feature** taken from top –j layer of encoder

- Transformer-based thresholdless multi-intent NLU framework w/ 3 multi-task learning
                                                                – Intent cls, IND cls\*, Sot cls
- The output of each upper j-layer in the encoder is used to <u>generate multi-grain representations</u> at different levels of granularity (passed through FC and just sum them up)

# Related Works ; Baseline

**Legend:**
- intent 0
- PlayMusic
- AddToPlaylist
- intent 1
- intent 2
- intent 3

$\hat{y} = 1$

① Softmax ② Sigmoid

*Let* **k = 1**,

*Choose* **top − 1 intents**

[CLS]    $T_1$    $T_2$   …   $T_N$

## Pretrained Language Model (**PLM**)
( e.g. BERT, RoBERTa, ELECTRA, …)

[CLS]    $Tok_1$    $Tok_2$   …   $Tok_N$

**Tokenizer**

Since it is conversational speech,

omissions, contractions, and ungrammatical structures may occur.

play harry styles falling and add it to favorite playlist

① **Intent Number Detection**
- **Input**: 1 utterance
- **Output**: 1 label (# of intent)

② **Multi-intent classification**
- **Input**: 1 utterance
- **Output**: 1 or more label(s)

**IND** : an auxiliary task, model detects the number of intents in each utterance

HYU 한양대학교 HANYANG UNIVERSITY

# Related Works ; Baseline

| Model | MixATIS | | | MixSNIPS | | |
|---|---|---|---|---|---|---|
| | Slot (F1) | Intent (Acc) | Overall (Acc) | Slot (F1) | Intent (Acc) | Overall (Acc) |
| SF-ID (*concat*) (2019) | 87.4 | 66.2 | 34.9 | 90.6 | 95.0 | 59.9 |
| Stack-Propagation (*thresh* = 0.5) (2019) | 87.8 | 72.1 | 40.1 | 94.2 | 96.0 | 72.9 |
| Joint Multiple ID-SF (*thresh* = 0.5) (2019) | 84.6 | 73.4 | 36.1 | 90.6 | 95.1 | 62.9 |
| AGIF (*thresh* = 0.5)(2020) | 86.7 | 74.4 | 40.8 | 94.2 | 95.1 | 74.2 |
| GL-GIN (*thresh* = 0.5)(2021) | **88.3** | 76.3 | 43.5 | 94.9 | 95.6 | 75.4 |
| SDJN (*thresh* = 0.5)(2022a) | 88.2 | 77.1 | 44.6 | 94.4 | 96.5 | 75.7 |
| SDJN+BERT (*thresh* = 0.5)(2022a) | 87.5 | 78.0 | 46.3 | 95.4 | 96.7 | 79.3 |
| Bert-baseline (*thresh* = 0.3) | 83.1 | 74.8 | 42.6 | 95.5 | 95.7 | 80.2 |
| Bert-baseline (*thresh* = 0.5) | 86.3 | 74.5 | 44.8 | 95.5 | 95.6 | 80.1 |
| Bert-baseline (*thresh* = 0.8) | 85.6 | 75.8 | 43.5 | 95.2 | 96.7 | 80.6 |
| **TFMN** (Bert-base) | 88.0 | **79.8** | **50.2** | **96.4** | **97.7** | **84.7** |

- **Transformer-Based** Models in MID: Achieving Unprecedented High Performance.
- Competitive Edge: Latest Models Contending in Subtle Decimal Point Differences.
- MID and Slot Filling: **Jointly-learning** (multi-task learning) in Advanced Research Fields.

# Related Works ; MixSNIPS & MixATIS

- **MixSNIPS** (Introduced by Qin et al. 2020)

  - Advanced SNIPS for multi-intent classification
  - Just concatenate sentences using "and" with different intents
    - Ratio of sentences: [1, 2, 3] intents [0.3, 0.5, 0.2]
  - Size: 50,000 utterances
  - Label: 7 intents (up to 3-label)

- \* Leaderboard
  - Intent Detection: 97.7 (accuracy)
  - Slot Filling: 96.4 (F1)

- **MixATIS** (Introduced by Qin et al. 2020)

  - Advanced ATIS for multi-intent classification
  - Just concatenate sentences using "and" with different intents
    - Ratio of sentences: [1, 2, 3] intents [0.3, 0.5, 0.2]
  - Size: 20,000 utterances
  - Label: 22 intents (up to 3-label)

- \* Leaderboard
  - Intent Detection: 76.3 (accuracy)
  - Slot Filling: 88.3 (F1)
  - Semantic Frame Parsing: 43.5 (accuracy)

sample ▼

**BookRestaurant**

```
89   book O
90   a O
91   reservation O
92   for O
93   my B-party_size_description
94   mommy I-party_size_description
95   and I-party_size_description
96   i I-party_size_description
97   at O
98   a O
99   restaurant B-restaurant_type
100  in O
101  central B-country
102  african I-country
103  republic I-country
104  and O
```

**PlayMusic**

```
105  then O
106  play O
107  the O
108  newest B-sort
109  melody B-music_item
110  on O
111  last B-service
112  fm I-service
113  by O
114  eddie B-artist
115  vinson I-artist
116  BookRestaurant#PlayMusic
```

\* intents (indicator: '#')

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- ## Overview

| | |
|---|---|
| **선정 이유** | − Multi-intent utterance 데이터를 구축하려는 시도 조사<br>: 기존 MixATIS, MixSNIPS에 대한 한계를 인지하고, 보다 현실적인 setting을 제안한 데이터셋 조사 |
| **특징** | − Multi-step human-annotated dataset<br>− 중국에서 구축, 중국어 사용 현실을 우선적으로 반영<br>− 후속 쿼리(follow-up query) 생성시 crowd-worker가 의도적으로 유관한 내용을 생성하도록 유도한 점이 데이터 품질을 향상시켰을 것으로 추측 |
| **프로젝트 기여 가능성** | − DialogUSR: initial query에 연결되는 follow-up query를 human-annotation → queries를 모두 merge한 utterance에 대한 데이터셋 제안<br>   − Intent detection 대상 utterance를 한 문장으로 제한하지 않고, **여러 문장의 연속**으로 처리<br>   − 후속 쿼리 생성 시 **자연스럽게 주제 전환이 발생**된다는 등, 데이터 구축 시 참고 가능한 일부 통계 확인<br>− 중국어 setting을 우선하여 구축된 데이터 → 영어 또는 한국어 등 **다른 언어 현실에 general하게 반영될 수 있는 내용인지 확인 불가**<br>− Slot filling에 대한 labeling에 대해서는 future work으로 제안<br>− 의미적으로 동시 발생 가능성이 높아 현실적인 multi-intent utterance가 구축됨<br>− 기존 setting보다는 다양한 접속사를 활용<br>   − DialogUSR의 inference: multi-intent 발화를 single-intent로 분리 → SID 진행<br>     (현업에서 쉽고 간편한 확장성을 목적으로 함) → split이 쉽도록 query들이 연결된 경향이 있음<br>   − 이 논문에서는 데이터셋 제안을 위한 구축 과정 설명 및 split 성능에 대해서만 보고됨 (MID 성능에 대해서는 보고하지 않음) |

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- ## Summary

  - Multi-intent Detection Dataset (**w/o slot-filling**) contains 11.6k high quality instances that cover 23 domains
  - Single-intent user input → Multi-intent user input (reformulation)
  - Multi-step human-annotation (🔽 samples)

| Domain | Query1 | Query2 | Query3 | Query4 | Query |
|--------|--------|--------|--------|--------|-------|
| Attraction | Recommend fun and cheap places near Kunshan | Where is Kunshan Miaofeng Pagoda | Navigate to Kunshan Miaofeng Pagoda | None | Recommend fun and cheap places near Kunshan. In addition to this, where is the Kunshan Miaofeng Pagoda? Navigate to Kunshan Miaofeng Pagoda. |
| TV | Hunan Satellite TV entertainment program | Is there any funny variety show on Hunan Satellite TV | Is there any program hosted by he Jiong on Hunan Satellite TV | Does Hunan Satellite TV have any funny movies at 8:00 pm | Hunan Satellite TV entertainment program. Are there any funny variety shows? Is there any program hosted by he Jiong? Are there any funny movies at 8 pm? |
| Train | Tomorrow's train to Zhengzhou | What is the latest train number to Zhengzhou | How much is the train fare to Zhengzhou | How many trains are there a day to Zhengzhou? | The train to Zhengzhou tomorrow. I would also like to know what time is the latest train? And how much is the fare? How many trains are there a day? |
| Weather | What is the weather in Linyi | How is the air quality in Linyi | What is the temperature in Linyi in the next week | Will it rain in Linyi | What is the weather in Linyi? In addition, I would like to know how is the air quality in Linyi? And what about the temperature in the coming week? Besides this, will it rain? |
| Restaurant | Search for nearby western restaurants | Which is the closest western restaurant to me | Which western restaurant has the highest score | None | Search for nearby western restaurants. Also I would like to know which one is closest to me? Which has the highest rating. |
| Flight | Shanghai to Beijing flights | Buy me the earliest flight from Shanghai to Beijing | Buy me the cheapest flight from Shanghai to Beijing | How many flights are there from Shanghai to Beijing every day | Flights from Shanghai to Beijing. Buy me the earliest flight. Buy me the cheapest ticket. How many flights are there in a day? |

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- **Suggestion**

  - Dataset Construction

  1. **Initial Query Collection**
     - Based on Single-intent datasets: SMPECDT2, RiSAWOZ3
     - Except queries that have excessive length or too verbose/repetitive in terms of semantics

  2. **Follow-up Query Creation**
     - Ask for "imagine they are eliciting multiple intents in a single complex user query"
     - Instruct annotators to write ~3 subsequent queries on what they need or would like to know about according to the initial query
     - 37.3% multi-intent queries involve topic switching
       * <u>conforms</u> to the user behavior in the real-world

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- **Suggestion**
  - Dataset Construction
  1. **Initial Query Collection**
  2. **Follow-up Query Creation**
  3. **Query Aggregation**
     - Findings of Pilot Study
       - Lack of variations in the conjunctions btw sub-queries: 'and', 'or', 'then', 'also', …
       - Lack of diversity & naturalness of the derived query in the human-only annotation
     - Connect sub-queries
       - Concatenate 2 consecutive queries w/ or w/o <u>text filler</u> (50% chance)
         - pre-define templates: "first of all", "I also would like to know", "finally", …
       - Even being locally coherent, the derived multi-intent query may still exhibit some global incoherence and syntactic issues
     - Post-check the sentence fluency of aggregated queries by GPT-2 (117M) model

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- **Suggestion**
  - Dataset Construction
  1. **Initial Query Collection**
  2. **Follow-up Query Creation**
  3. **Query Aggregation**
  4. **Query Completion**
     - 62.5% follow-up queries occur coreferences(2.4%) and omissions
       * do NOT explicitly ask the annotator to use incomplete queries in 2. follow-up query
     - recover missing information by annotators, correctly annotating 8 out of 10 cases

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- **Suggestion**

  - Dataset Construction

  - Annotation Settings

    - Be paid 0.6$ per datapoint, which is more than prevailing local minimum wage

    - split the entire annotation procedure into multiple rounds
      → hire another group judges to post-check the quality of annotated dataset
      → filter unqualified instances after each round

  - Statistics: Total 11,669 instances

    - Aggregated multi-intent complex query is assembled 3.6 single-intent and comprise 36.7 Chinese characters.

    - Initial query/1$^{st}$/2$^{nd}$/3$^{rd}$ follow query's length: 11.9/12.3/12.4/10.8

# Related Works

📄 **DialogUSR)** Meng et al. DialogUSR: Complex Dialogue Utterance Splitting and Reformulation for Multiple Intent Detection, EMNLP Findings, 2022

- **Suggestion**

  - Dataset Construction

  - Annotation Settings

  - Task Overview
    - ○ Q1: Multi-intent query in DialogUSR
    - ○ Q2: Split Multi-intent query to Single-intent queries     sequence generation task
    - ○ Q3: Delete Conjunctions / Q4: Recover Missing Info. (coreference, omission)
    - ○ Q5: Recover 1st split-query
    - ○ Q6: Recover 2nd split-query and concatenate w/Q5    independent!
    - ○ Q7: Recover 3rd split-query and concatenate w/Q6

  - End-to-end Generative Models: Q1 → Q4

  - 2-stage Generative Models
    - (once)　2-stage model: Q1 → Q2 → Q4
    - (casual) 2-stage model: Q1 → Q2 → [Q5 → Q6 → Q7]

almost 100% split

| Model | MixSNIPS | | MixATIS | |
|---|---|---|---|---|
| | BLEU | EM | BLEU | EM |
| T5-base | 99.46 | 95.13 | 96.94 | 74.88 |
| T5-large | 99.60 | 97.64 | 98.52 | 88.77 |
| T5-xl | **99.62** | **98.14** | **99.87** | **98.55** |

Input(Q1): 查询周五下午厦门到南京的动车需要多长时间 然后查一下那边的特色美食
Check the high-speed train from Xiamen to Nanjing on Friday afternoon, how long does the journey take, then check out the special food there.

Split(Q2): 查询周五下午厦门到南京的动车 [SP] 需要多长时间 [SP] 然后查一下那边的特色美食
Check the high-speed train from Xiamen to Nanjing on Friday afternoon [SP] how long does the journey take [SP] then check out the special food there.

Delete(Q3): 查询周五下午厦门到南京的动车 [SP] 需要多长时间 [SP] 查一下那边的特色美食
Translation is the same as above

Complete(Q4): 查询周五下午厦门到南京的动车 [SP] 厦门到南京的动车需要多长时间 [SP] 查一下南京的特色美食
Check the high-speed train from Xiamen to Nanjing on Friday afternoon [SP] How long does it take to travel from Xiamen to Nanjing in high-speed train [SP] Check out the special cuisine in Nanjing

Causal Complete:
Step1(Q5): 查询周五下午厦门到南京的动车 => 查询周五下午厦门到南京的动车
Check the high-speed train from Xiamen to Nanjing on Friday afternoon => Translation is the same as above

Step2(Q6): 查询周五下午厦门到南京的动车 [SP] 需要多长时间 => 厦门到南京的动车需要多长时间
Check the high-speed train from Xiamen to Nanjing on Friday afternoon [SP] how long does the journey take => How long does it take to travel from Xiamen to Nanjing in high-speed train

Step3(Q7): 查询周五下午厦门到南京的动车 [SP] 需要多长时间 [SP] 然后查一下那边的特色美食 => 查一下南京的特色美食
Check the high-speed train from Xiamen to Nanjing on Friday afternoon [SP] how long does the journey take [SP] then check out the special food there => Check out the special cuisine in Nanjing

End-to-end (E2E): Q1 → Q4
Two-stage (Once): Q1 → Q2 → Q4
Two-stage (Casual): Q1 → Q2 → [Q5 → Q6 → Q7]

HYU 한양대학교 HANYANG UNIVERSITY

# Benchmark Datasets Analysis (1/2)

Multiple Single-intent utterances

⌄

1 Multi-intent utterance

**ATIS**

give me the round trip flights
from cleveland to miami next wednesday

atis_flight

**+**

give me the fares for round trip flights
from cleveland to miami next wednesday

atis_airfare

⌄⌄

**MixATIS**

give me the round trip flights
from cleveland to miami next wednesday
and give me the fares for round trip flights
from cleveland to miami next wednesday

Libo Qin, Xiao Xu et al. AGIF: An Adaptive Graph-Interactive Framework for Joint Multiple Intent Detection and Slot Filling. EMNLP 2020, Findings.

The dataset relies on only a few specific connectors ("and", "and then", "and also")
when concatenating 2 or more single-intent utterances.

→ Real-world conversations often involve **more varied and complex ways of combining intents**

# Benchmark Datasets Analysis (2/2)

Multiple Single-intent utterances

1 Multi-intent utterance

ATIS

give me the round trip flights
from cleveland to miami next wednesday

atis_flight

**+**

give me the fares for round trip flights
from cleveland to miami next wednesday

atis_airfare

MixATIS

Libo Qin, Xiao Xu et al. AGIF: An Adaptive Graph-Interactive Framework for Joint Multiple Intent Detection and Slot Filling. EMNLP 2020, Findings.

give me the round trip flights
from cleveland to miami next wednesday
and give me the fares for round trip flights
from cleveland to miami next wednesday

give me the fares and round trip flights
from cleveland to miami next wednesday

We are focused on constructing our own set that better mirrors natural language usage
to provide more **challenging** and **realistic** resources
for training and testing multi-intent detection models.

한양대학교
HANYANG UNIVERSITY

22

# Discussion 1

# Concatenation: Single- to Multi- intent utterance

# Utterance Selection

HYU 한양대학교
HANYANG UNIVERSITY

# Categories of Concatenation Complexity-side

- **Complexity side**
  - Explicit Concatenation: use connectors during concatenation
    - → AND variants / Various Conjunctions
  - Implicit Concatenation: do NOT use connectors during concatenation
    - → Inherent Ambiguity / Gerund Phrase / Omission / Coreference

## Complexity side

**Explicit Concatenation**
- Only use connectors to concatenate utterances
  – and, and then, and also
  – before, additionally, …

**AND variants** — and, and then, and also

**Various Conjunctions** — ,(comma), ;(semicolon), or, before, after, additionally, finally

**Implicit Concatenation**
- Concatenate utterances without connectors
  – Inherent ambiguity
  – Omissions
  – Coreferences
  – Gerund phrase

**Inherent Ambiguity** — just joined two sentences directly without using a connector

**Gerund Phrase** — provide additional meaning by adding context or details to the action described by the gerund

**Omission** — to avoid redundancy, intentionally leave out repeated words for conciseness

**Coreference** — use different words or phrases to refer back to the same entities, ensuring coherence in the text

Natural Language Processing Lab., Hanyang University.

# Categories of Concatenation Methodology-side

- **Methodology side**
  - Manual Concatenation: rule-based concatenation approach
    - → AND variants / Various Conjunctions / Inherent Ambiguity / Gerund Phrase
  - Generative Concatenation: concatenation by using generative language model
    - → Omission / Coreference

**Methodology side**

*and, and then, and also* → **AND variants**

*, (comma), ;(semicolon),*
*or, before, after, additionally, finally* → **Various Conjunctions**

just joined two sentences directly
without using a connector → **Inherent Ambiguity**

provide additional meaning by adding context or details to the
action described by the gerund → **Gerund Phrase**

to avoid redundancy,
intentionally leave out repeated words for conciseness → **Omission**

use different words or phrases to refer back to the same entities,
ensuring coherence in the text → **Coreference**

| Naïve Approach |
| --- |
| • Approach employed by MixX |

| Manual Approach |
| --- |
| • Rule-based concatenation technique |

| Generative Approach |
| --- |
| • Concatenate utterances by using generative language model, ChatGPT |

# **Categories of Concatenation** Methodology-side

- **Example for each concatenation approach**

| | utterance1 | utterance2 | intent1 | intent2 | Concatenation result |
|---|---|---|---|---|---|
| **Naïve** | i want to put this song in my new boots playlist | what films are going to be playing at harkins theatres at zero a m | AddToPlaylist | SearchScreeningEvent | i want to put this song in my new boots playlist and what films are going to be playing at harkins theatres at zero am |
| **Manual** | please show me all airports in denver | can you list costs of denver rental cars | atis_airport | atis_ground_fare | please show me all airports in denver listing costs of denver rental cars |
| **Generative** | play some theme songs from 1974 | play the movie white christmas | PlayMusic | SearchCreativeWork | play some theme songs from 1974 and the movie white Christmas |
| | clear my to do list | repeat my to do list | todo_list_update | todo_list | i need to clear my to-do list and then repeat it |

- Connectors excluded from Manual Concatenation

If ~
as ~ / but~
to + verb
…
→ Eliminate intent

I'd like to ~
Search for ~
Could you ~
…
→ paraphrasing

# Categories of Concatenation Methodology-side

- ## Complexity side

  - Explicit Concatenation: use connectors during concatenation
  - Implicit Concatenation: do NOT use connectors during concatenation

- ## Methodology side

  - Manual Concatenation: rule-based concatenation approach
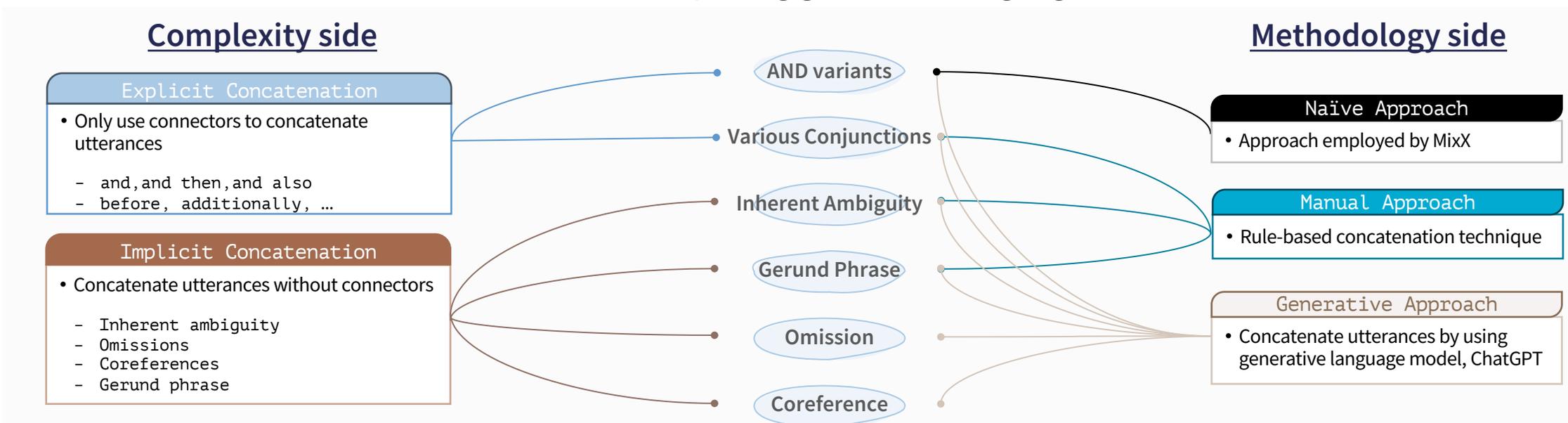  - Generative Concatenation: concatenation by using generative language model

# (Intuitive) ChatGPT for Concatenation (1/2)

- ## Prompt Engineering for ChatGPT Concatenation

```
You are a native English speaker.
[Task Definition] Combine 2 or 3 sentences as one single sentence.
[Goal] The focus is on creating a single sentence that captures the essence of both ideas without unnecessary redundancy.
[Instructions] - Avoid adding just punctuation.
               - Don't paraphrase.
               - Don't compromise the meaning of each sentence.
               - Don't capitalize all characters.
               - Don't replace numbers with radix.
               - Maintain the intent of each sentence.
               - Don't forget that if a sentence starts with a verb, it's a statement.
               - Do NOT use conjunctions like 'and'
               - Don't print '[Good Answer]'
               - Don't print intent directly.
[Example 1]
        My dog is playful (dog's feature) + My dog loves chasing balls (dog's feature)
        [Good Answer] My playful dog loves chasing balls
        [Bad Answer] My dog is playful, and my dog loves chasing balls
        [Bad Answer] My dog is playful, and also loves chasing balls.
[Example 2]
        They finished the project(project done) + They had time(taking time)
        [Good Answer] With time on their hands, they finished the project
        [Bad Answer] They finished the project, and they had time
        [Bad Answer] They finished the project and had time
[Example 3]
        She answered the phone (answering)) + She was making dinner (cooking)
        [Good Answer] While answering the phone, she was making dinner
        [Bad Answer] She answered the phone, and She was making dinner
Combine the following sentences naturally. Inside the parentheses is the intent of each sentence. :
{utt1} (intent: {intents[0]})  + {utt2} (intent: {intents[1]})
```

Returning results that don't follow the explicit constraints we gave ChatGPT

Few-shot setting

# (Intuitive) ChatGPT for Concatenation (2/2)

- **Failure of Using ChatGPT**
  - Over 90% of the 100 randomly sampled instances were connected by 'and'
  - Failed to maintain intent in at least 15 instances
    - Intent removed / changed
    - Fail to merge the utterances

| | Utterance 1 | Utterance 2 | Concatenation | Implication |
|---|---|---|---|---|
| failed | weather for frankfort (`GetWeather`) | 3 out of 6 for the last album (`RateBook`) | frankfort's weather gets a 3 out of 6 rating for the last album | Intent removed |
| | what day of the week do flights from nashville to tacoma fly on (`atis_day_name`) | flight number from houston to dallas (`atis_flight_no`) | flights from nashville to tacoma fly on what day of the week and what is the flight number from houston to dallas | Intent changed |
| | what is mci (`atis_abbreviation`) | list la (`atis_city`) | combine the sentences: "what is mci?" and "list la". | Failed to merge |
| succeeded | play the playlist funtime activity (`AddToPlaylist`) | add adele onto funtime activity playlist (`PlayMusic`) | play the playlist, adding adele to the funtime activity playlist | Gerund phrase |
| | how do i freeze my bank account (`account_blocked`) | do you know why my bank account is frozen (`freeze_account`) | how do i freeze my bank account, and do you know why it is frozen | Coreference |

Table 1: Failed (Top) and succeeded (Bottom) results of the **Generative Approach** and their implications.

Only few cherry-picked examples were, brilliantly, implicitly concatenated, which is what we intended.

# Overview of 2 Methods to Utterance Concatenation

**1** For Manual Approach

**2** Utterance Selection

Randomly select 2~3 utterances

**Manual Concatenation**

- AND variants
- Various Conjunctions
- Inherent Ambiguity
- Gerund Phrase

**Naïve Approach**
- Approach employed by MixX

**Manual Approach**
- Rule-based concatenation technique

**1** Single-Intent Datasets

- Banking77
- ATIS
- CLINC150
- SNIPS

**Preprocessing**
- Lowercase
- Remove punctuation

**2** For Generative Approach

**3** Utterance Selection

Select 2~3 similar utterances based on cosine similarity

**Generative Concatenation**

- AND variants
- Various Conjunctions
- Inherent Ambiguity
- Gerund Phrase
- Omission
- Coreference

ChatGPT Bard Llama

**Generative Approach**
- Concatenate utterances by using generative language model, ChatGPT

**Filtering Concatenated Utterances**

Enhance the quality of outcomes from generative LLMs by implementing defined criteria and incorporating reviews from human experts to filter the results.

1 ✓ Develop metrics to evaluate integration diversity and complexity
✓ collect results above a set threshold.

2 ✓ Collect only the incorrect predictions made by the best model in multi-intent detection.

3 ✓ Human-expert review of the generated output filtered by the above two steps.

**4** BlendX

- BlendBanking77
- BlendATIS
- BlendCLINC150
- BlendSNIPS

**Construct multi-intent utterance dataset**

- To complete the final dataset, we have multi-intent utterances
- concatenated by manual approach
- and those generated via a generative approach, refined with human filtering.

---

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose 2 approaches :

1 **Manual Approach:** Concatenate utterances without using connectors, or if necessary, employ a various range of options.

2 **Generative Approach:** Explore methods to extend **ChatGPT**'s capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.

# Overview of 2 Methods to Utterance Concatenation



**1**

### Single-Intent Datasets

Banking77   ATIS

CLINC150   SNIPS

### Preprocessing

- Lowercase
- Remove punctuation

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose 2 approaches :

① **Manual Approach:** Concatenate utterances without using connectors, or if necessary, employ a various range of options.

② **Generative Approach:** Explore methods to extend **ChatGPT**'s capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.

# Overview of 2 Methods to Utterance Concatenation

**2**

**1** For Manual Approach

**Single-Intent Datasets**

Banking77    ATIS
CLINC150    SNIPS

**Preprocessing**

- Lowercase
- Remove punctuation

**Utterance Selection**

$Y$

$X$

$Z$

- Randomly select 2~3 utterances

**Manual Approach**

AND variants

Various Conjunctions

Inherent Ambiguity

Gerund Phrase

**Naïve Approach**
- Approach employed by MixX

**Manual Approach**
- Rule-based concatenation technique

Select 2~3 similar utterances based on cosine similarity

Coreference

**3**

✓ Human-expert review of the generated output filtered by the above two steps.

**BlendX**

endBanking77    BlendATIS
endCLINC150    BlendSNIPS

**Construct multi-intent utterance dataset**

To complete the final dataset, we have multi-intent utterances concatenated by manual approach and those generated via a generative approach, refined with human filtering.

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose 2 approaches :

**1** **Manual Approach:** Concatenate utterances without using connectors, or if necessary, employ a various range of options.

**2** **Generative Approach:** Explore methods to extend ChatGPT's capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.
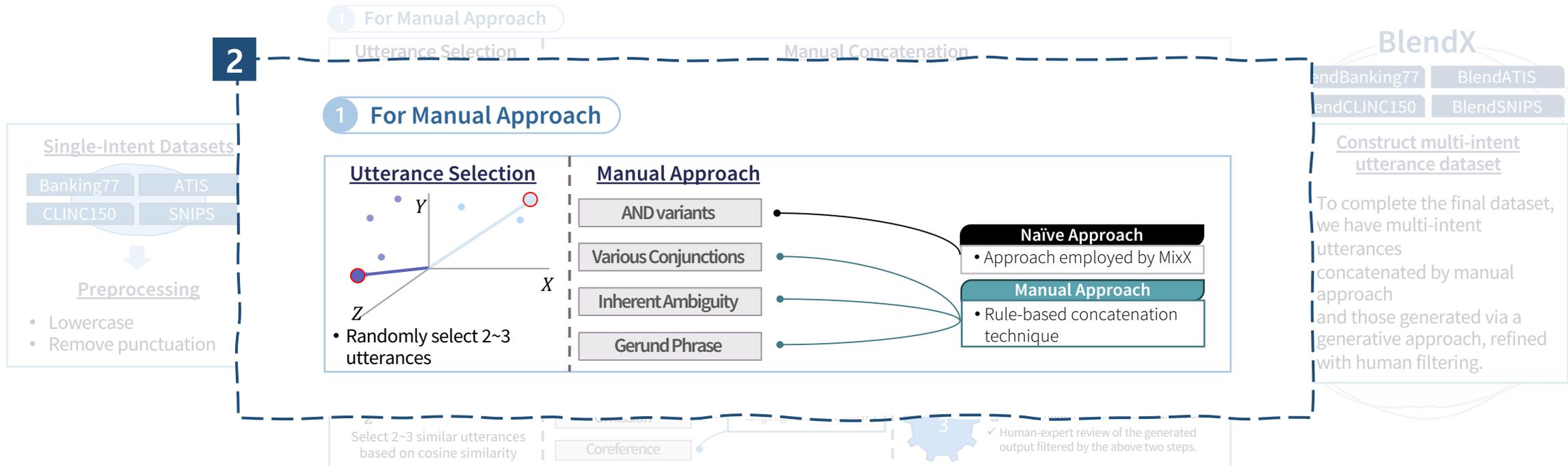
# Overview of 2 Methods to Utterance Concatenation

**2**



**① For Manual Approach**

**Utterance Selection** | **Manual Approach**

| and, and then, and also | AND variants |
| ‘,’, ‘;’, or, before, after, … | Various Conjunctions |
| Remove connector | Inherent Ambiguity |
| Use ‘-ing’ form | Gerund Phrase |

Select 2~3 similar utterances based on cosine similarity

give me the round trip flights
from cleveland to miami next wednesday

**+**

give me the fares for round trip flights
from cleveland to miami next wednesday

give me the round trip flights
from cleveland to miami next wednesday
~~and~~ give me the fares for round trip flights
from cleveland to miami next wednesday

---

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose **2 approaches** :

**①** **Manual Approach:** Concatenate utterances without using connectors, or if necessary, employ a various range of options.

**②** Generative Approach: Explore methods to extend ChatGPT's capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.
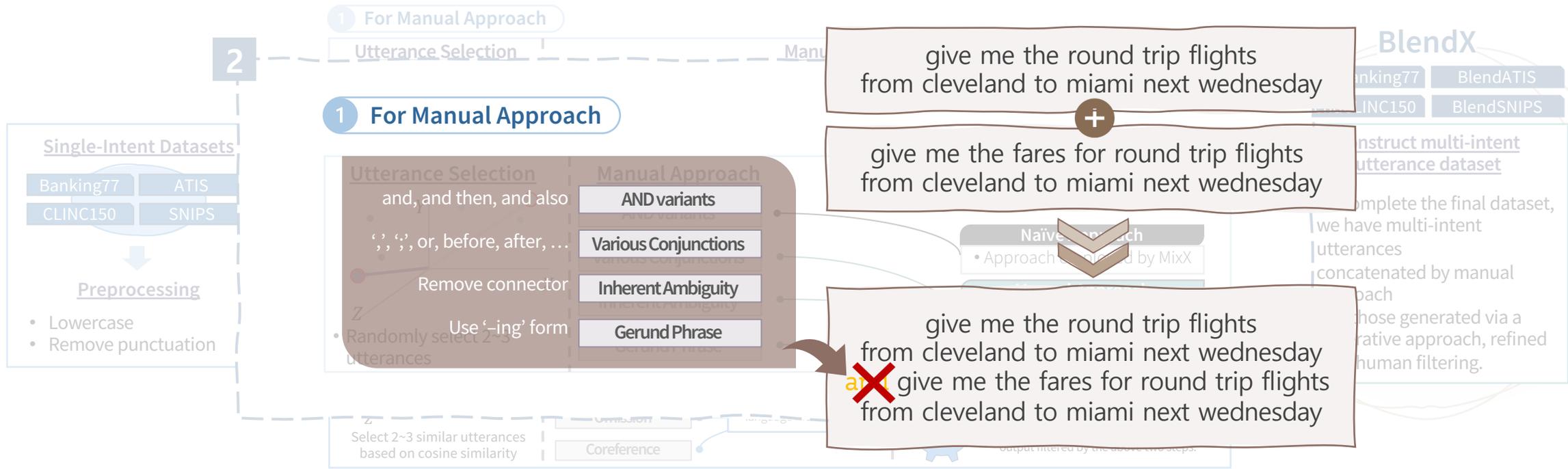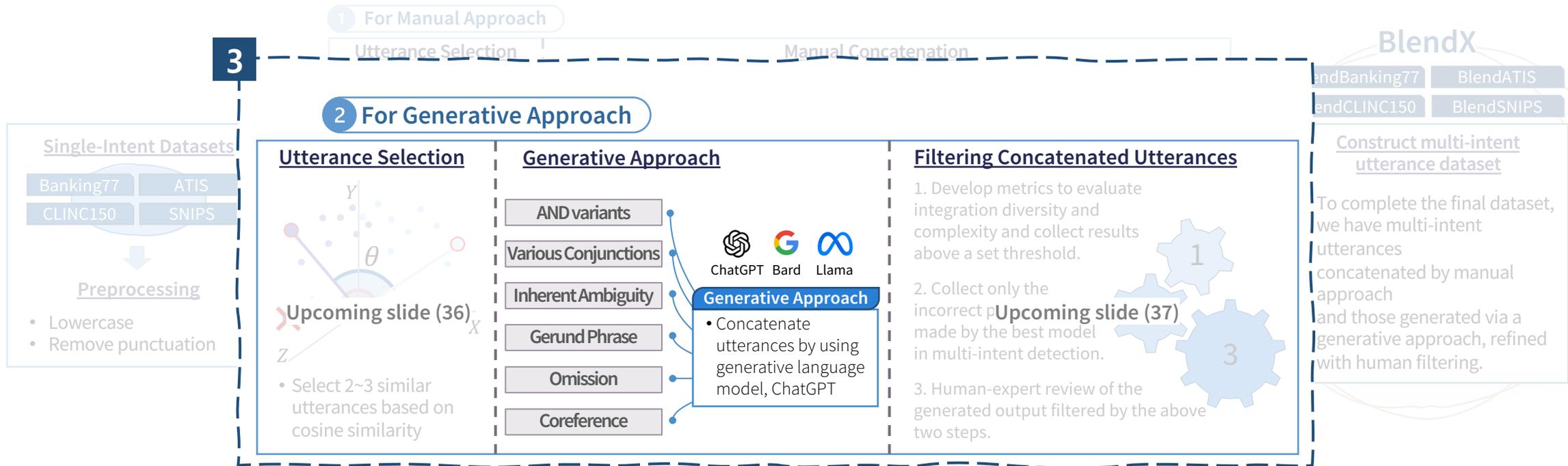
# Overview of 2 Methods to Utterance Concatenation

**3**

① For Manual Approach

Utterance Selection

Manual Concatenation

**BlendX**

BlendBanking77 | BlendATIS
BlendCLINC150 | BlendSNIPS

② **For Generative Approach**

**Single-Intent Datasets**

Banking77 | ATIS
CLINC150 | SNIPS

**Preprocessing**

- Lowercase
- Remove punctuation

**Utterance Selection**

$Y$

$\theta$

$X$

$Z$

**Upcoming slide (36)**

- Select 2~3 similar utterances based on cosine similarity

**Generative Approach**

| AND variants |
| Various Conjunctions |
| Inherent Ambiguity |
| Gerund Phrase |
| Omission |
| Coreference |

ChatGPT  Bard  Llama

**Generative Approach**
- Concatenate utterances by using generative language model, ChatGPT

**Filtering Concatenated Utterances**

1. Develop metrics to evaluate integration diversity and complexity and collect results above a set threshold.

2. Collect only the incorrect p**Upcoming slide (37)** made by the best model in multi-intent detection.

3. Human-expert review of the generated output filtered by the above two steps.

1

3

**Construct multi-intent utterance dataset**

To complete the final dataset, we have multi-intent utterances concatenated by manual approach and those generated via a generative approach, refined with human filtering.

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose 2 approaches :

① Manual Approach: Concatenate utterances without using connectors, or if necessary, employ a various range of options.

② **Generative Approach:** Explore methods to extend **ChatGPT**'s capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.

HYU 한양대학교 HANYANG UNIVERSITY

# Overview of 2 Methods to Utterance Concatenation

**3**

**①** For Manual Approach

Utterance Selection

Manu...

BlendX

Banking77 | BlendATIS

CLINC150 | BlendSNIPS

**②** **For Generative Approach**

**Single-Intent Datasets**

Banking77 | ATIS

CLINC150 | SNIPS

**Preprocessing**

- Lowercase
- Remove punctuation

**Utterance Selection**

**Generative Approach**

**①** **For Manual Approach**

ChatGPT also use these
techniques for concatenation

leave out repeated words

use different words or phrases to
refer back to the same entities

AND variants

Various Conjunctions

Inherent Ambiguity

Gerund Phrase

Omission

Coreference

ChatGPT  Bard  Llama

ChatGPT

give me the round trip flights
from cleveland to miami next wednesday

**+**

give me the fares for round trip flights
from cleveland to miami next wednesday

give me the fares and round trip flights
from cleveland to miami next wednesday

construct multi-intent
utterance dataset

complete the final dataset,
we have multi-intent
utterances
concatenated by manual
approach
and those generated via a
generative approach, refined
human filtering.

Without generating new multi-intent utterances and ensuring they fit within the existing intent space,
we propose 2 approaches :

**①** Manual Approach: Concatenate utterances without using connectors, or if necessary, employ a various
range of options.

**②** Generative Approach: Explore methods to extend **ChatGPT**'s capabilities for producing more coherent
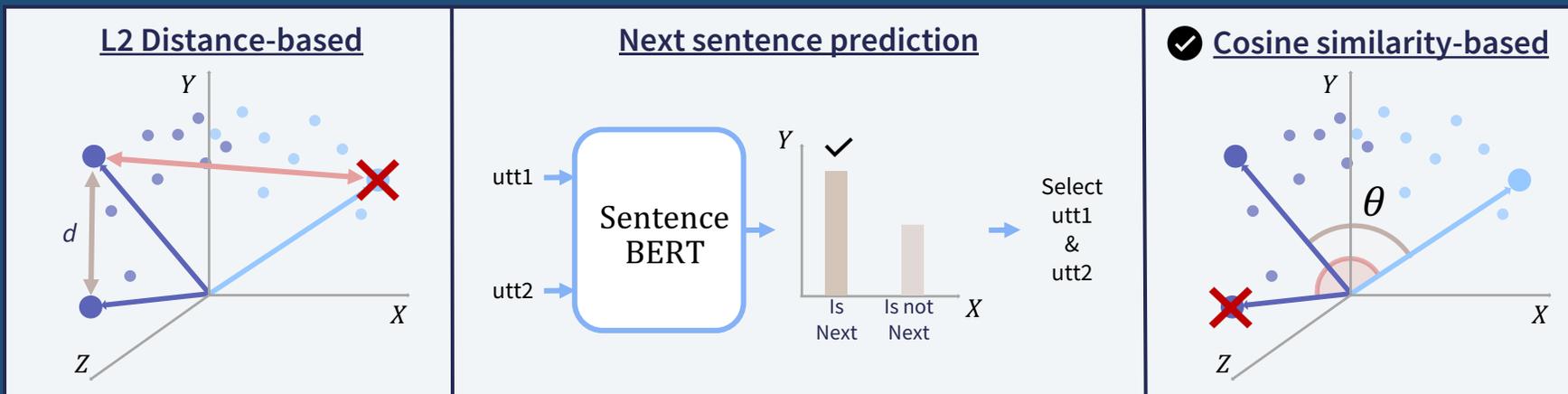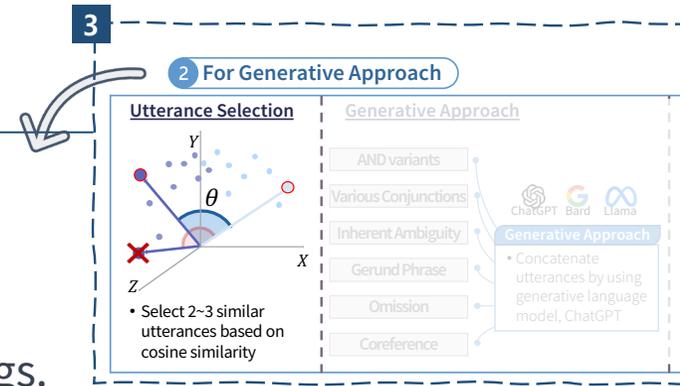multi-intent utterances by concatenating 2 or more single-intent utterances.

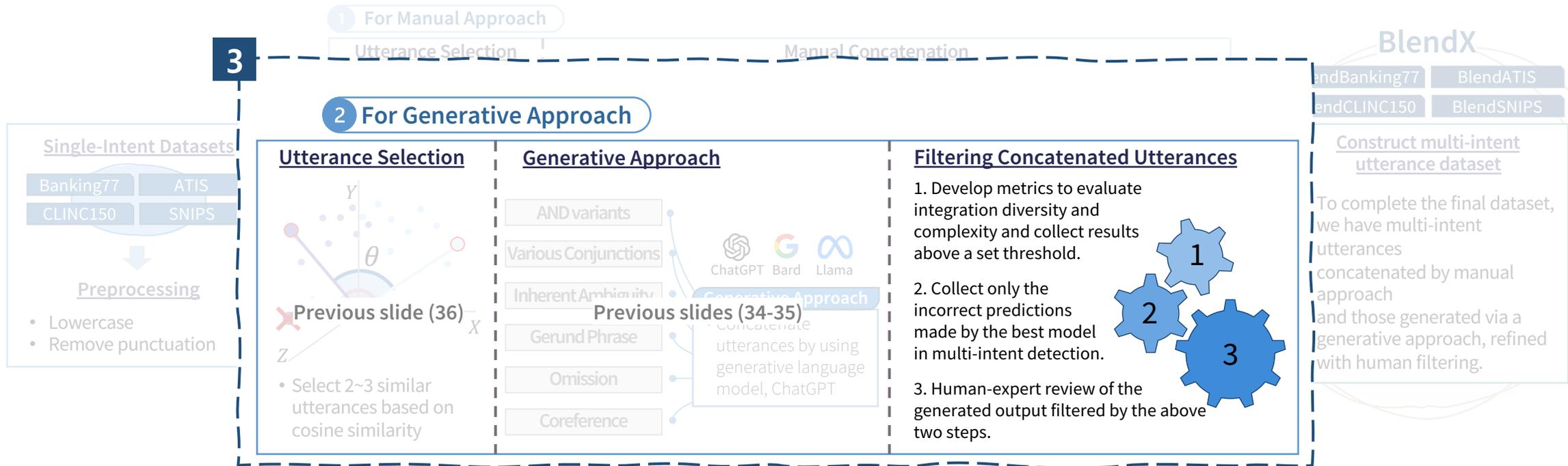# Utterance Selection for ❷ Generative Approach

## • Process

1. Generate embeddings for each single-intent utterance using SentenceBERT.

2. Select utterances for concatenation based on high similarity between embeddings.
   * Chosen utterances will have different intents.

## • Selection approach

- L2 Distance-based: Select utterances with close proximity in embedding space.

- Next sentence prediction: Binary classification of whether a given pair of utterances are sequential.

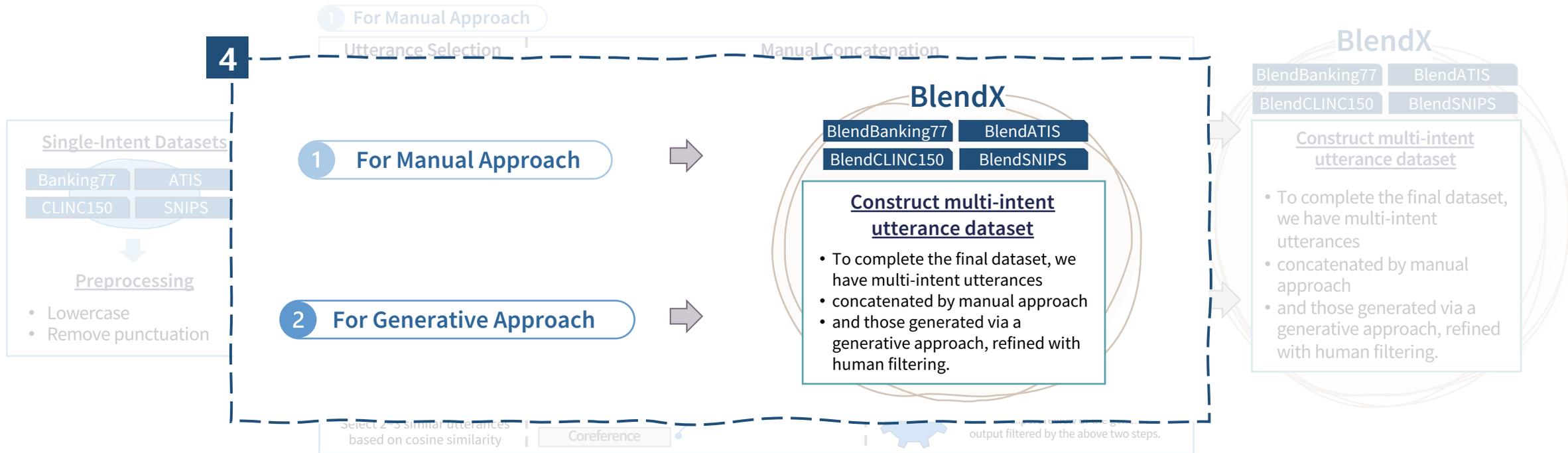✔ - Cosine similarity-based: Choose utterances with high cosine similarity between embeddings.



**L2 Distance-based**        **Next sentence prediction**        ✔ **Cosine similarity-based**

# Overview of 2 Methods to Utterance Concatenation

**3**

① For Manual Approach

Utterance Selection    Manual Concatenation

② For Generative Approach

**BlendX**

BlendBanking77    BlendATIS
BlendCLINC150    BlendSNIPS

**Single-Intent Datasets**

Banking77    ATIS
CLINC150    SNIPS

**Preprocessing**

- Lowercase
- Remove punctuation

**Utterance Selection**

$Y$

$\theta$

$X$

$Z$

**Previous slide (36)**

- Select 2~3 similar utterances based on cosine similarity

**Generative Approach**

AND variants

Various Conjunctions

Inherent Ambiguity

Gerund Phrase

Omission

Coreference

ChatGPT    Bard    Llama

**Generative Approach**

**Previous slides (34-35)**

- Concatenate utterances by using generative language model, ChatGPT

**Filtering Concatenated Utterances**

1. Develop metrics to evaluate integration diversity and complexity and collect results above a set threshold.

2. Collect only the incorrect predictions made by the best model in multi-intent detection.

3. Human-expert review of the generated output filtered by the above two steps.

1
2
3

**Construct multi-intent utterance dataset**

To complete the final dataset, we have multi-intent utterances concatenated by manual approach and those generated via a generative approach, refined with human filtering.

Without generating new multi-intent utterances and ensuring they fit within the existing intent space, we propose 2 approaches :

① Manual Approach: Concatenate utterances without using connectors, or if necessary, employ a various range of options.

② Generative Approach: Explore methods to extend **ChatGPT**'s capabilities for producing more coherent multi-intent utterances by concatenating 2 or more single-intent utterances.

HYU 한양대학교 HANYANG UNIVERSITY

# Overview of 2 Methods to Utterance Concatenation

**4**

① For Manual Approach

Utterance Selection                          Manual Concatenation

**Single-Intent Datasets**

Banking77    ATIS
CLINC150    SNIPS

**Preprocessing**
- Lowercase
- Remove punctuation

**①  For Manual Approach**

**②  For Generative Approach**

## BlendX

BlendBanking77    BlendATIS
BlendCLINC150    BlendSNIPS

**Construct multi-intent utterance dataset**
- To complete the final dataset, we have multi-intent utterances
- concatenated by manual approach
- and those generated via a generative approach, refined with human filtering.

## BlendX

BlendBanking77    BlendATIS
BlendCLINC150    BlendSNIPS

**Construct multi-intent utterance dataset**
- To complete the final dataset, we have multi-intent utterances
- concatenated by manual approach
- and those generated via a generative approach, refined with human filtering.

Select 2~3 similar utterances based on cosine similarity   Coreference   output filtered by the above two steps.

---

## *BlendX* : Complex multi-intent detection with blended patterns

| Dataset | # of intents | Training | Dev | Test | Total |
|---|---|---|---|---|---|
| BlendSNIPS | 7 | 50,625 | 2,613 | 2,615 | 55,853 |
| BlendATIS | 18 | 20,250 | 1,125 | 1,125 | 22,500 |
| BlendBanking77 | 77 | 36,390 | 2,009 | 2,021 | 40,420 |
| BlendCLINC150 | 147 | 54,899 | 2,889 | 2,977 | 60,765 |

∑ (total) = 179,538

- Source Dataset : SNIPS, ATIS, Banking77, CLINC150
- Random selection for **Manual** Concatenation Approach
- Cosine Similarity-based selection for **Generative** Concatenation Approach

# Discussion 2

# 1) 3 Custom Metrics

# 2) Comparison BlendX to MixX with SOTA baselines

# 3) Verify semantics differences in vector space

HYU 한양대학교
HANYANG UNIVERSITY

# Evaluation #1 – 3 Custom Metrics (1/4)

- ## 3 Custom Metrics

  - $utt$: concatenated utterance with 2 or more intents

  - $n$: Number of single-intent utterances used for concatenation

### $W(utt, n)$
#### Word count

$$\mathrm{W}(utt, n) \overset{\text{def}}{=} \mathbf{1}_{\mathbb{Z}-\mathbb{N}}\left(|utt|_{word} - \sum_{i=1}^{n} |utt_i|_{word}\right).$$

Check if the **word count** difference
before and after
an utterance concatenation is
zero or negative

( to ascertain a decrease in word count )

### $C(utt, n)$
#### Conjunction

$$\mathrm{C}(utt, n) \overset{\text{def}}{=} \mathbf{1}_{\mathbb{Z}-\mathbb{N}}\left(|utt|_{conj} - \sum_{i=1}^{n} |utt_i|_{conj}\right).$$

Verify if the number of **conjunctions**
before and after
an utterance changes to zero or less

( to determine the elimination
or reduction of conjunctions )

* **conjunctions** such as:
and, or, before, after,
additionally, finally, ',', ';'

### $P(utt, n)$
#### Pronoun

$$\mathrm{P}(utt, n) \overset{\text{def}}{=} \mathbf{1}_{\mathbb{N}}\left(|utt|_{pron} - \sum_{i=1}^{n} |utt_i|_{pron}\right).$$

Assess if the difference in **pronoun**
**count** before and after
an utterance is one or more

( to identify the usage of pronouns )

* **pronoun** such as :
it, them, their, theirs, this, that,
those, these

An implicitly concatenated utterance is likely to receive 1 in the metrics evaluation.

# Evaluation #1 – 3 Custom Metrics (2/4)

- **Example of applicating 3 metrics**

| | Concatenation | utt1 | utt2 | Difference | Metric |
|---|---|---|---|---|---|
| **example #1** | add another song to my 88 keys playlist playing it | play my 88 keys playlist | add another song to my 88 keys playlist | | |
| **Words** | 10 | 5 | 8 | 10 - (5 + 8) = -3 | $W(\cdot, 2) = 1$ |
| **Conjunctions** | 0 | 0 | 0 | 0 - (0 + 0) = 0 | $C(\cdot, 2) = 1$ |
| **Pronouns** | 1 | 0 | 0 | 1 - (0 + 0) = 1 | $P(\cdot, 2) = 1$ |
| **example #2** | i need to clear my to-do list and then repeat it | clear my to do list | repeat my to do list | | |
| **Words** | 11 | 5 | 5 | 11 - (5 + 5) = 1 | $W(\cdot, 2) = 0$ |
| **Conjunctions** | 1 | 0 | 0 | 1 - (0 + 0) = 1 | $C(\cdot, 2) = 0$ |
| **Pronouns** | 1 | 0 | 0 | 1 - (0 + 0) = 1 | $P(\cdot, 2) = 1$ |

| **Utterance 1** | play my 88 keys playlist (`PlayMusic`) | | | |
|---|---|---|---|---|
| **Utterance 2** | add another song to my 88 keys playlist (`AddToPlaylist`) | | | |

| **Strategies** | **Concatenation Results** | $W(utt, 2)$ | $C(utt, 2)$ | $P(utt, 2)$ |
|---|---|---|---|---|
| **Explicit Concatenation** | play my 88 keys playlist **and** also add another song to my 88 keys playlist | 0 | 0 | 0 |
| **Implicit Concatenation** | | | | |
| Inherent Ambiguity | play my 88 keys playlist add another song to my 88 keys playlist | 1 | 1 | 0 |
| Omissions | play my 88 keys playlist and add another song | 1 | 0 | 0 |
| Coreferences | play my 88 keys playlist and add another song to it | 1 | 0 | 1 |
| Gerund Phrase | add another song to my 88 keys playlist playing it | 1 | 1 | 1 |

Table 3: Various concatenation classes, accompanied by their examples and respective metric values.

# Evaluation #1 – 3 Custom Metrics (3/4)

- **Results using 3 metrics for each approach** (Naïve, Manual, Generative)

*MixX*

| Metric | SNIPS | | | ATIS | | | Banking77 | | | CLINC150 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Naïve | Manual | Generative | Naïve | Manual | Generative | Naïve | Manual | Generative | Naïve | Manual | Generative |
| $W(utt, 2)(\uparrow)$ | 0% | **37%** | 29% | 0% | **36%** | 18% | 0% | **46%** | 37% | 0% | **48%** | 28% $\uparrow$ |
| $C(utt, 2)(\uparrow)$ | 0% | **56%** | 10% | 0% | **52%** | 15% | 0% | **50%** | 27% | 0% | **56%** | 32% $\uparrow$ |
| $P(utt, 2)(\uparrow)$ | 0% | 0% | **7%** | 0% | 0% | **8%** | 0% | 0% | **13%** | 0% | 0% | **6%** |

Table 4: Comparative analysis of the three concatenation approaches: Naïve, Manual, and Generative. Notably, the Manual method demonstrates pronounced efficiency in reducing utterance length.

Our approach, incorporating both **manual** and **generative** methods, achieves
a more diverse range of explicit and implicit concatenation compared to existing techniques.

- Notably, **MixX** did not involve **implicit** concatenation. ☐
  "Naïve" refers to the original construction method of **MixX**, meaning concatenation using only `and`, `and then`, and `and also`.
- Particularly, **manual concatenation** often resulted in shorter utterance lengths. ☐
- Conversely, **generative concatenation** uniquely led to the use of pronouns. ☐

# Evaluation #1 – 3 Custom Metrics (4/4)

**3**

**2** For Generative Approach

Utterance Selection



• Select 2~3 similar utterances based on cosine similarity

- **Results using 3 metrics for <u>Generative</u> approach w/ <u>utterance selection</u>**

| Metric | SNIPS | | ATIS | | Banking77 | | CLINC150 | |
|---|---|---|---|---|---|---|---|---|
| | Random | Sim. | Random | Sim. | Random | Sim. | Random | Sim. |
| Cosine sim. | 0.105 | 0.746 | 0.214 | 0.758 | 0.212 | 0.748 | 0.093 | 0.749 |
| Error rate (↓) | 16% → | **14%** | 41% → | **10%** | 22% → | **9%** | 19% → | **13%** |
| $W(utt, 2)(\uparrow)$ | 27.38% | **44.87%** | 10.17% | **27.78%** | **34.62%** | 30.77% | 30.86% | **31.03%** |
| $C(utt, 2)(\uparrow)$ | **8.33%** | 1.28% | 3.39% | **4.44%** | **28.21%** | 15.38% | **25.93%** | 3.45% |
| $P(utt, 2)(\uparrow)$ | 3.57% | **10.26%** | 1.69% | **12.22%** | 10.26% | **20.88%** | 3.70% | **14.94%** |

ChatGPT Concatenation Failure Rate

Table 2: Comparison of Random and Similarity-Based (Sim.) utterance selection across datasets when applied to ChatGPT. We find that Sim. leads to a reduced error rate in ChatGPT's data generation.

- Compared to random selection, similarity-based selection
  - **reduces the error rate** by anywhere from 2% to as much as 31%. ☐
  - **increases the use of pronouns** ☐ and **decreases the word count** in most cases.
- Increased frequency of implicit concatenation, especially omission or coreference, which naturally leads to **increased use of simple conjunction 'and' variants** to ensure semantic clarity. ☐
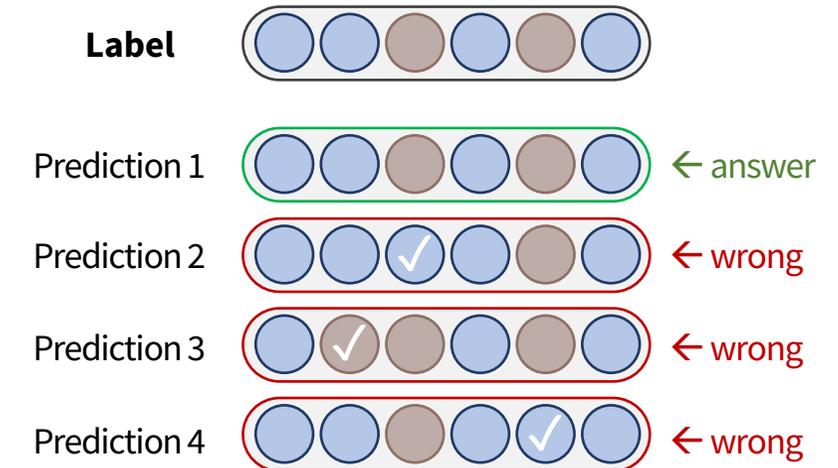
# Evaluation #2 – BlendX vs. MixX (1/2)

- **Comparison BlendX to MixX (SOTA baselines)**

| Model | Option | | Dataset (Metric: accuracy) | | | |
|---|---|---|---|---|---|---|
| | Training | Test | SNIPS | ATIS | Banking77 | CLINC150 |
| TFMN | MixX | MixX | 95.68* ±0.57 | 77.98* ±0.57 | 76.61 ±1.17 | 85.88 ±1.03 |
| | MixX | BlendX | 52.51 ±1.86 | 42.51 ±1.48 | 37.31 ±0.81 | 42.45 ±2.40 |
| | BlendX | BlendX | 94.93 ±0.85 | 76.50 ±0.83 | 63.99 ±0.81 | 77.96 ±0.82 |
| SLIM | MixX | MixX | 95.97* ±0.23 | 77.10* ±0.28 | 83.71 ±0.88 | 88.67 ±0.56 |
| | MixX | BlendX | 93.51 ±0.18 | 72.80 ±1.48 | 69.89 ±0.46 | 73.39 ±2.46 |
| | BlendX | BlendX | 95.73 ±0.86 | 76.92 ±0.84 | 75.30 ±0.71 | 85.62 ±0.51 |
| gpt-3.5-turbo | - | MixX | 81.68 | 40.30 | 30.90 | 49.22 |
| | - | BlendX | 76.18 | 38.84 | 22.67 | 37.55 |

- **Accuracy** in Multi-label Classification
  : only considered it correct in cases of a **exact match**.



For various SOTA models, we consistently observe a huge performance drop ↘ on our **BlendX** datasets with explicit as well as implicit concatenations.

- 3-Baseline: implemented without slot-filling part
  - ✓ **TFMN** : predict # of intents $k$, and then top-$k$ intents over the probability distribution
  - ✓ **SLIM** : threshold-based classification model using sigmoid function
  - ✓ **ChatGPT** : OpenAI's generative model ($gpt-3.5-turbo-0613$)

# Evaluation #2 – BlendX vs. MixX (2/2)

- ## ChatGPT ICL prompt
  - prompts as simple and easy to understand as possible.

```
You are an Intent Detection Model on single utterance.

[Task Definition]
    Detect single or more intent(s) of each utterance, but you can only classify UP TO 3
      most plausible intents on 1 utterance.
[Labels] atis_airport, atis_ground_service, atis_abbreviation, atis_city, atis_aircraft,
      atis_ground_fare, atis_flight, atis_airfare, atis_meal, atis_distance, atis_cheapest,
      atis_capacity, atis_restriction, atis_quantity, atis_airline, atis_flight_no,
      atis_flight_time, atis_day_name
[Answer format]
    If more than one, concatenate with '#', such as {Label}#{Label}.
      e.g. atis_ground_fare#atis_distance


[Example 1]
      [Utterance] what is restriction ap80
      [Answer] atis_restriction
[Example 2]
      [Utterance] what does the fare code qx mean , what is the distance between Pittsburgh
              airport and downtown pittsburgh and what is restriction ap80
      [Answer] atis_abbreviation#atis_distance#atis_restriction

Detect a single or up to 3 intent(s) on this following utterance: utt
```
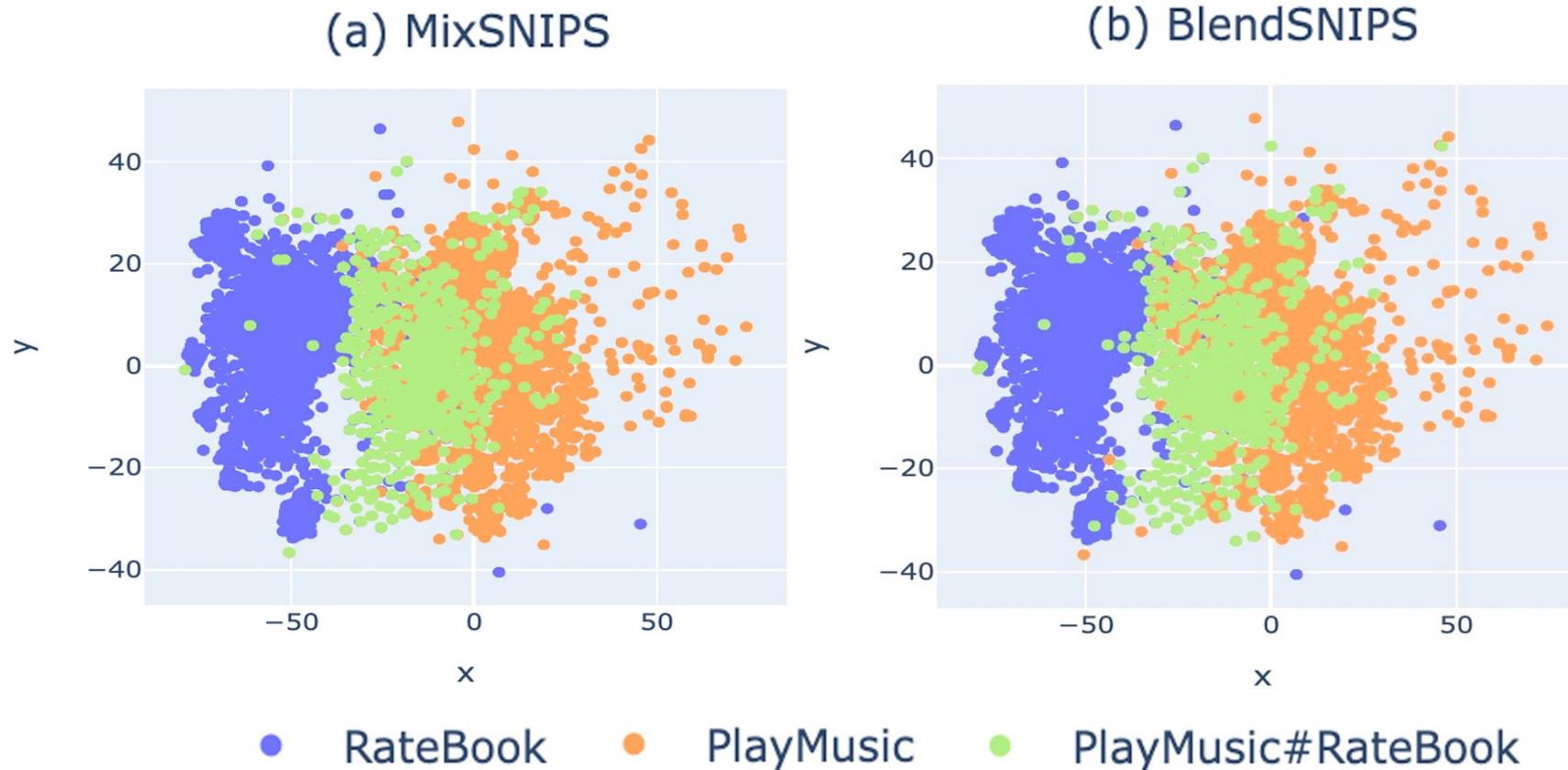
If ChatGPT returned results that did not follow explicit constraints (maximum 3 intents, answer format, etc.), we post-processed and measured performance.

한양대학교
HANYANG UNIVERSITY

# Evaluation #3 – Visualization

- **Visualization of <u>MixX</u> and <u>BlendX</u> utterances on 2-dimensional space**
  - BlendX's concatenated utterances preserve the semantics of both source utterances.

# Conclusion

# Main Findings

# Limitation

# Future Work

HYU 한양대학교
HANYANG UNIVERSITY

# Main Findings

- **Identified limitations in existing multi-intent datasets**

  - MixX: Reliance on explicit concatenation through the `'and'` connector.

- **BlendX:** Constructing a more complex and realistic multi-intent dataset

  - Proposed 3 novel concatenation approaches
                          : Naïve, Manual, Generative

  - Beyond random sentence selection,
    applied a similarity-based strategy
    in the **generative** concatenation approach.

  - Designed 3 statistical metrics for comparing and
    validating **BlendX** against the existing **MixX**: W, C, P

  - Upcoming dataset release
    : Extensions of **MixX** (**CLINC150**/**Banking77**) and
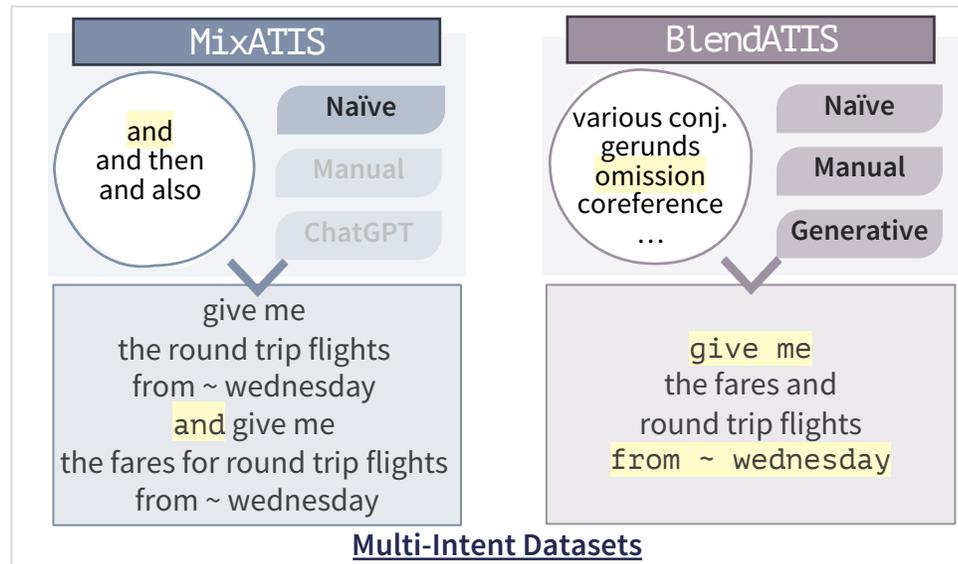    new publication of the **BlendX** dataset.

**#1 Selection**

Single-Intent Datasets

Banking77  ATIS✓
CLINC150   SNIPS

give me the round trip flights
from ~ wednesday          atis_flight

give me the fares for round trip flights
from ~ wednesday          atis_airfare

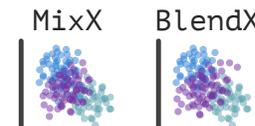**#2 Concatenation**

MixATIS

and
and then
and also

Naïve
Manual
ChatGPT

give me
the round trip flights
from ~ wednesday
and give me
the fares for round trip flights
from ~ wednesday

BlendATIS

various conj.
gerunds
omission
coreference
…

Naïve
Manual
Generative

give me
the fares and
round trip flights
from ~ wednesday

**Multi-Intent Datasets**

**#3 Evaluation**

|          | Mix | Blend |
|----------|-----|-------|
| $W(utt)$ | 0   | 1     |
| $C(utt)$ | 0   | 0     |
| $P(utt)$ | 0   | 0     |

**custom metric**

ChatGPT

**baseline evaluation**

MixX   BlendX

**visualization**

# Limitations

- **3 Custom Metrics**
  - Semantic Complexity is **not** considered.
    : Even though these metrics have improved, we cannot determine if this utterance is **semantically** complex.
    - e.g. i'd like to improve my credit score (*improve_credit_score*) + can you help me find my credit score (*credit_score*)
      - → i'd like to improve my credit score can you help me find it

  - Correlation between metrics is **not** considered.
    - If there is a pronoun in the concatenated utterance, does the word count decrease after concatenating?
    - If there are no conjunctions in the concatenated utterance, does the word count decrease after concatenating?
    - If the word count is reduced after concatenating, does the concatenated utterance have pronouns or no conjunctions?

- **Single dataset issue: Label overlapping**
  - CLINC150: *improve_credit_score, credit_score*
  - Banking77: *getting_virtual_card, get_disposable_virtual_card*

Cecilia Ying and Stephen Thomas. Label Errors in BANKING77.
ACL 2022 Workshop on Insights from Negative Results in NLP

HYU 한양대학교
HANYANG UNIVERSITY

# Limitations

- ## Limitation of Manual approach

| utt1 | utt2 | concatenation | implication |
|------|------|---------------|-------------|
| what is the least expensive fare from boston to salt lake city (atis_fare) | what are the fares for ground transportation in denver (atis_ground_fare) | what is the least expensive fare from boston to salt lake city **finally** what are the fares for ground transportation in denver | 두 의문문을 finally로 연결 |
| why isn't my id being verified (unable_to_verify_identity) | my top up was denied in the app (top_up_failed) | why isn't my id being verified **or** my top up was denied in the app | 의문문과 평서문을 or로 연결 |

- ## Limitation of Generative approach

| utt1 | utt2 | concatenation | implication |
|------|------|---------------|-------------|
| weather for Frankfort (GetWeather) | 3 out of 6 for the last album (RateBook) | frankfort's weather gets a 3 out of 6 rating for the last album | 'RateBook' 삭제 |
| what day of the week do flights from nashville to tacoma fly on (atis_day_name) | flight numbers from houston to dallas (atis_flight_no) | flights from nashwille to tacoma fly on what day of the week and what is the flight number from houstom to dallas | 'atis_day_name' → 'atis_flight' |
| what is mci (atis_abbreviation) | list la (atis_city) | combine the sentences: "what is mci" and "list la" | 결합 실패 |

# Future Work (ongoing, 1/3)

## 1. Multi-intent utterance Split → Intent Detection

End-to-end Generative Models: Q1 → Q4

2-stage Generative Models

    (once)  2-stage model
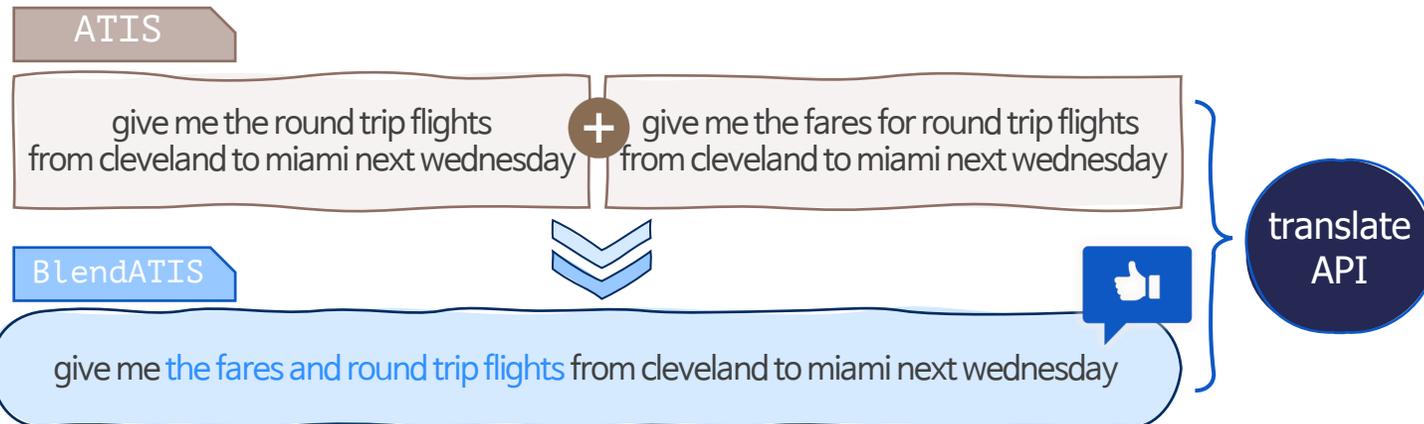
           : Q1 → Q2 → Q4

    (casual) 2-stage model     *almost 100% split*

           : Q1 → Q2 → [Q5 → Q6 → Q7]

| Model | MixSNIPS | | MixATIS | |
|---|---|---|---|---|
| | BLEU | EM | BLEU | EM |
| T5-base | 99.46 | 95.13 | 96.94 | 74.88 |
| T5-large | 99.60 | 97.64 | 98.52 | 88.77 |
| T5-xl | **99.62** | **98.14** | **99.87** | **98.55** |

## 2. Complex multi-intent utterance in Korean (w/ 지수, 정민, 정연)

ATIS

give me the round trip flights from cleveland to miami next wednesday **+** give me the fares for round trip flights from cleveland to miami next wednesday

translate API

BlendATIS

give me the fares and round trip flights from cleveland to miami next wednesday

---

1. **'~하고', '~하고 나서'와 같은 접속사(연결 어미)가 문장에 포함된 경우의 처리**

    - '에어컨 켜고 18도로 설정해줘' → '에어컨 켜고' + '18도로 설정해줘'

2. **Complex multi-intent 발화의 한국어 적용**

    - '밖에 미세먼지 좋으면 창문 다 열어줘'

Natural Language Processing Lab., Hanyang University.

HYU 한양대학교 HANYANG UNIVERSITY

# Future Work (ongoing, 2/3)

arXiv2024 / Do Large Language Model Understand Multi-Intent Spoken Language? / Yin et al.

## 3. ICL for Multi-intent Detection (w/ 정연, 영우)

- model: gpt-3.5-turbo-0613
- dataset: MixATIS
- Experiment (3-shot)
  - Naïve Prompt (from BlendX)
  - Role-assigned Prompt
  - Dynamic Few-shot Prompt
  - Multi-step

| Model | MixATIS | | | MixSNIPS | | |
|---|---|---|---|---|---|---|
| | Slot(F1) | Intent(Acc) | Overall(Acc) | Slot(F1) | Intent(Acc) | Overall(Acc) |
| Stack-Propagation (Qin et al., 2019) | 87.8 | 72.1 | 40.1 | 94.2 | 96.0 | 72.9 |
| AGIF (Qin et al., 2020b) | 86.9 | 72.2 | 39.2 | 93.8 | 95.1 | 72.7 |
| GL-GIN (Qin et al., 2021b) | 87.2 | 75.6 | 41.6 | 93.7 | 95.2 | 72.4 |
| SDJN (Chen et al., 2022) | 88.2 | 77.1 | 44.6 | 94.4 | 96.5 | 75.7 |
| CLID (Huang et al., 2022) | 88.2 | 77.5 | 49.0 | 94.3 | 96.6 | 75.0 |
| SSRAN (Cheng et al., 2023) | 89.4 | 77.9 | 48.9 | 95.8 | **98.4** | 77.5 |
| SDJN(BERT) | 87.5 | 78.0 | 46.3 | 95.4 | 96.7 | 79.3 |
| CLID(Roberta) | 85.9 | 80.5 | 49.4 | 96.0 | 97.0 | 82.2 |
| ChatGPT (5-shot) | 64.0 | 54.1 | 14.6 | 62.9 | 83.9 | 12.7 |
| Vicuna-7B-v1.5 (Peng et al., 2023) | 83.3 | 79.5 | 47.3 | 95.7 | 97.6 | 78.9 |
| Llama-2-7B-chat (Touvron et al., 2023) | 86.5 | 82.4 | 51.1 | 95.7 | 96.9 | 78.9 |
| Vicuna-13B-v1.5 (Peng et al., 2023) | 87.9 | **83.6***  | 50.8 | 95.9 | 97.5 | 80.7 |
| Llama-2-13B (Touvron et al., 2023) | 87.9 | 81.0 | 49.5 | **96.7** | 97.8 | **83.3*** |
| Mistral-7B-Instruct-v0.1 (Jiang et al., 2023) | **88.7** | 80.6 | **53.4*** | 95.6 | 97.6 | 79.8 |

| Experiment settings | | Accuracy |
|---|---|---|
| Naïve Prompt | Same as the prompt using in the experiments on BlendX paper | 33.60 |
| Role-assigned Prompt | According to OpenAI Official Document | 40.00 |
| Dynamic Few-shot Prompt | Use a {utterance, label} pair as a demonstration | 55.20 |
| Multi-step Prompt | Same as TFMN, but in-context learning setting | **79.90** |

# Future Work (ongoing, 2/3)

## 4. Improving BlendX (w/ 학부생 졸업프로젝트 #2)

- Thesis: N-gram 기반 similarity 측정을 통해 생략, 상호참조를 발생시키는 multi-intent 발화 데이터셋 구축

- Process
  1. Spacy 라이브러리를 사용하여 문장별 품사 및 구문분석 진행 → 명사/동사에 가중치 부여
  2. **1-gram 기반** 문장 별 유사도 계산
  3. 높은 순으로 pair 생성 (top-10)

• 유사도가 높고, 문장 합치기에 유리한 경우

| list all the airlines that fly into general mitchell international | list all the airlines that fly into general mitchell international | 0.9 |
| | list all the flights that arrive at general mitchell international | 0.6 |
| what is the earliest flight from boston to san francisco | what is the cheapest fare from boston to san francisco | 0.8 |

• 유사도가 높지만, 문장 합치기에 불리한 경우

| what flights are there from cleveland to miami on us air that arrive in miami before 4 pm | what round trip tickets are there from cleveland to miami on us air that arrive before 4 pm | 0.72 |
| | what is the cheapest first class fare from cleveland to miami on us air on february twenty fourth | 0.39 |

# Thank You

**Yejin Yoon**

HYU NLP Lab.
Hanyang University, South Korea

stillwithyou@hanyang.ac.kr

HYU 한양대학교
HANYANG UNIVERSITY