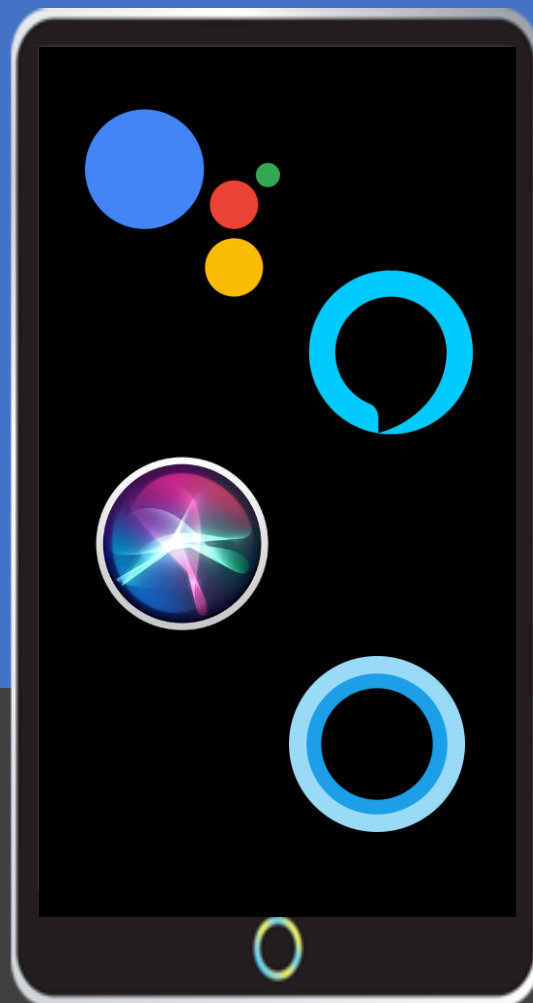


Robust and Scalable Conversational AI



?
=



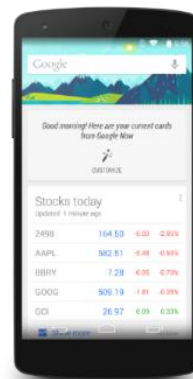
Yun-Nung (Vivian) Chen

Computer Science & Information Engineering
National Taiwan University

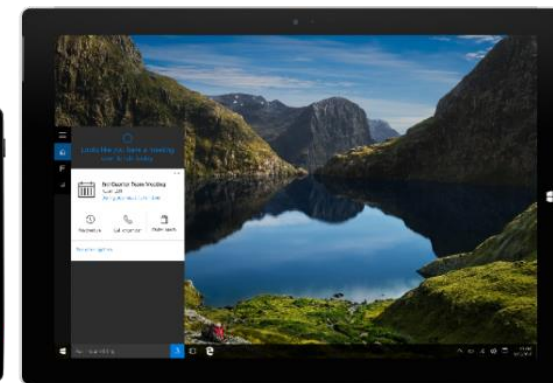
Language Empowering Intelligent Assistants



Apple Siri (2011)



Google Now (2012)
Google Assistant (2016)



Microsoft Cortana
(2014)



Amazon Alexa/Echo (2014)



Google Home (2016)

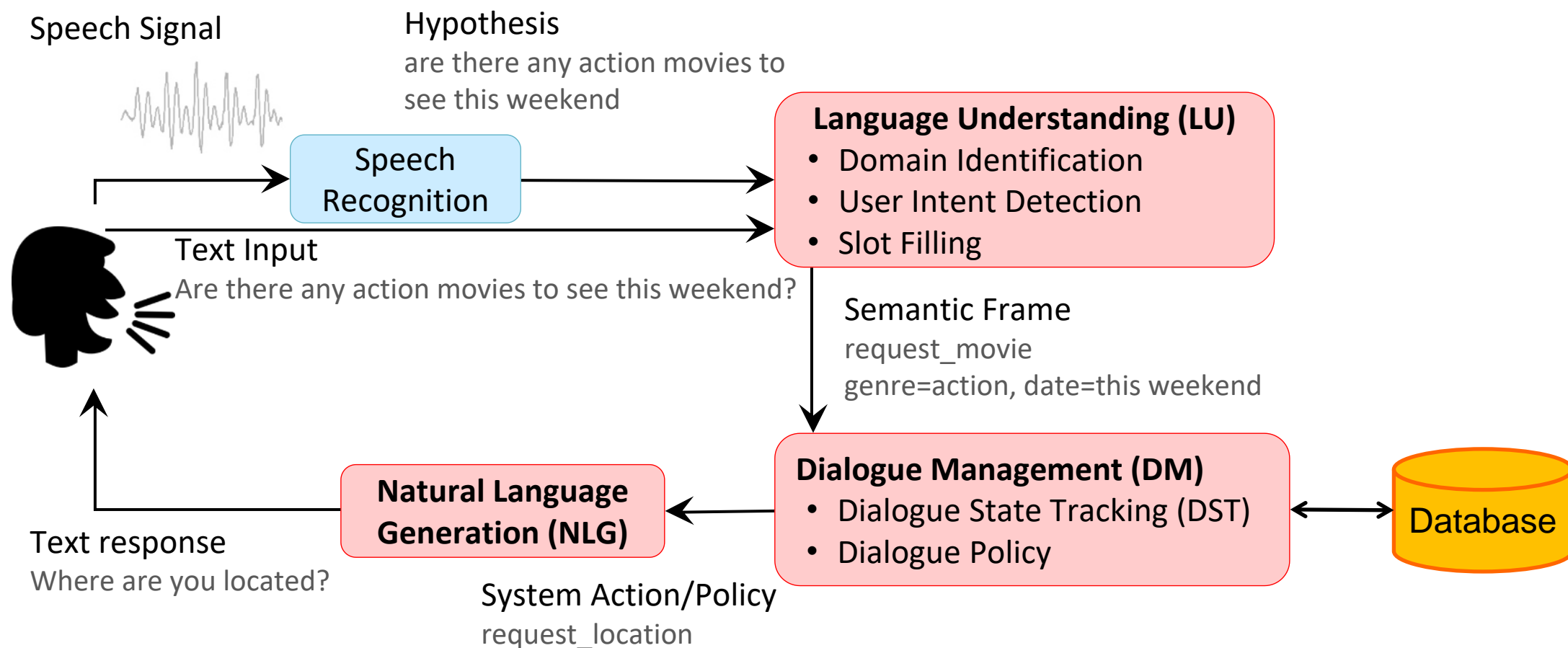


Apple HomePod (2017)



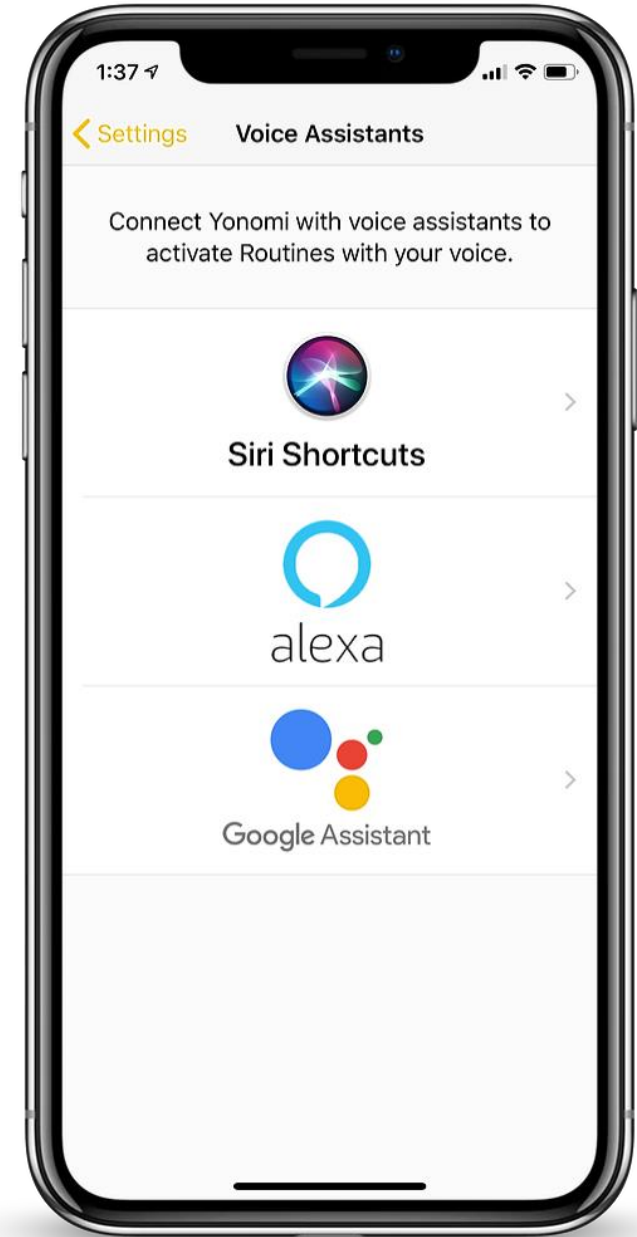
Facebook Portal (2019)

Task-Oriented Dialogue Systems ([Young, 2000](#))



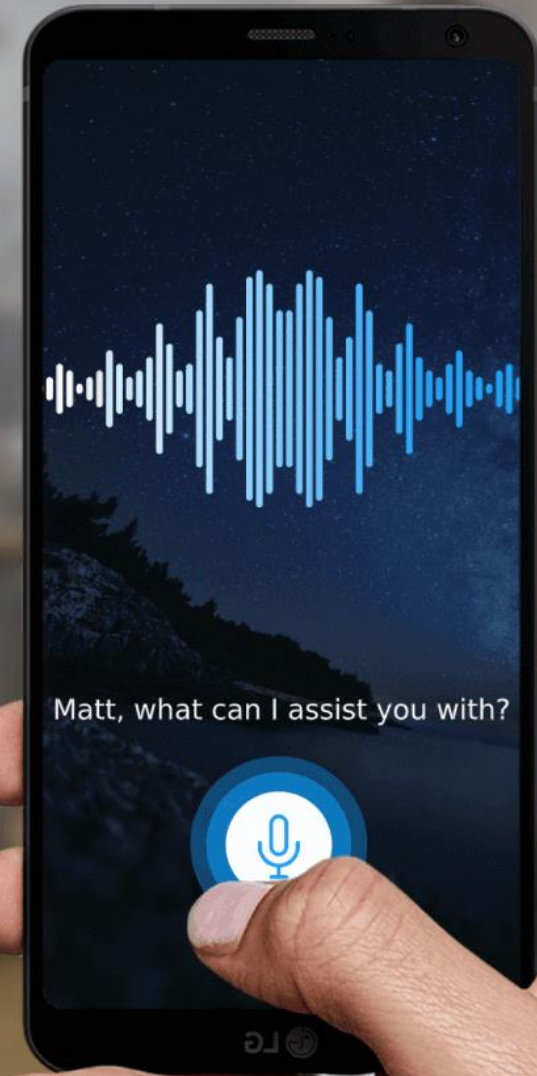
Recent Advances in NLP

- Contextual Embeddings (ELMo & BERT)
 - Boost many understanding performance with pre-trained language models

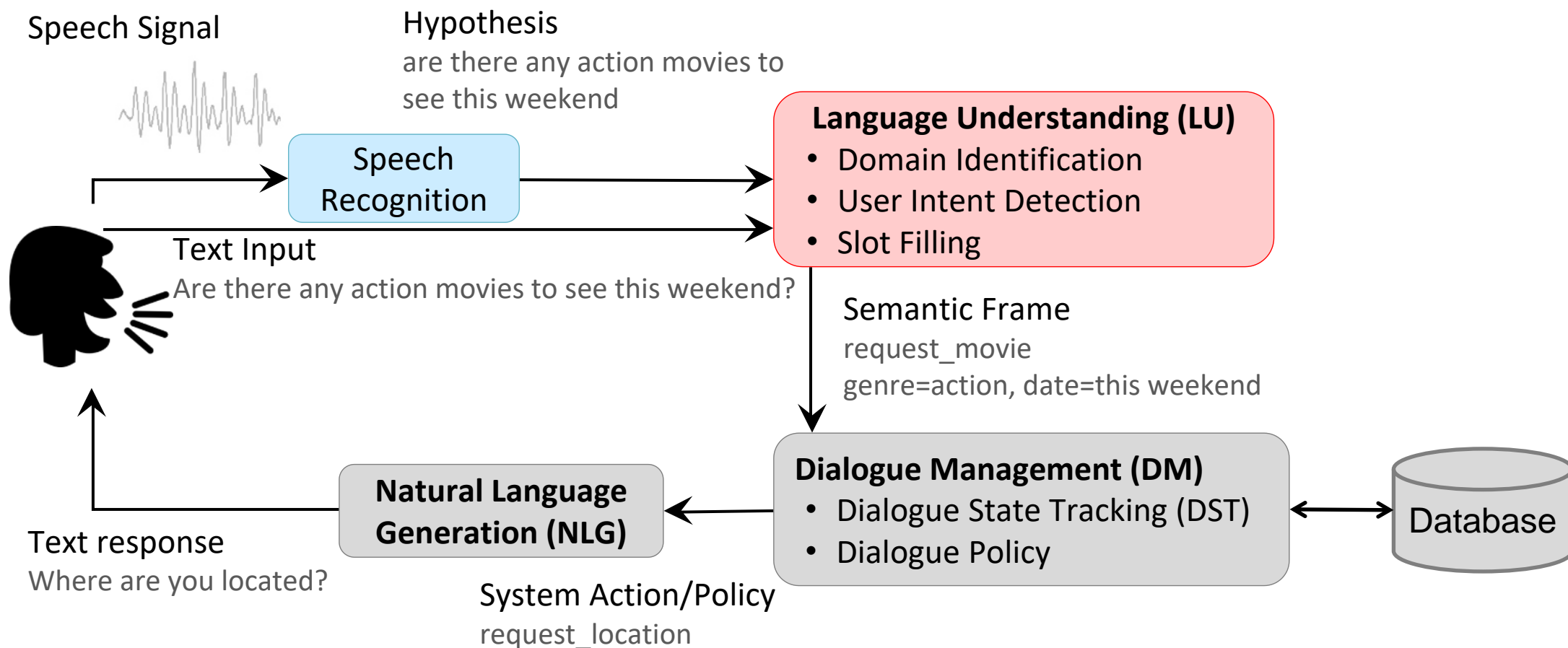




Lift all lights ~~X~~ to Morocco
List all flights tomorrow



Task-Oriented Dialogue Systems ([Young, 2000](#))



Mismatch between Written and Spoken Languages

Training

- Written language



Testing

- Spoken language
- Include recognition errors



- Goal: ASR-Robust Contextualized Embeddings
 - ✓ learning *spoken* contextualized word embeddings
 - ✓ better performance on *spoken* language understanding tasks

Solution: LatticeLM

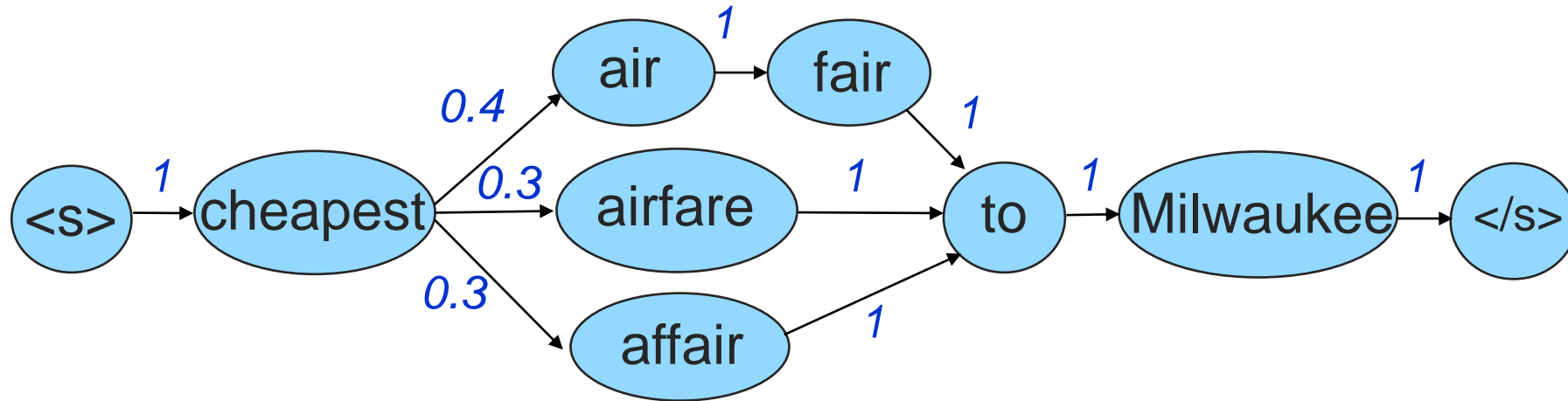
(Huang & Chen, ACL 2020)

9

Chao-Wei Huang and Yun-Nung Chen, “Learning Spoken Language Representations with Neural Lattice Language Modeling,” in *Proceedings of The 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, 2020.

ASR Lattices for Preserving Uncertainty

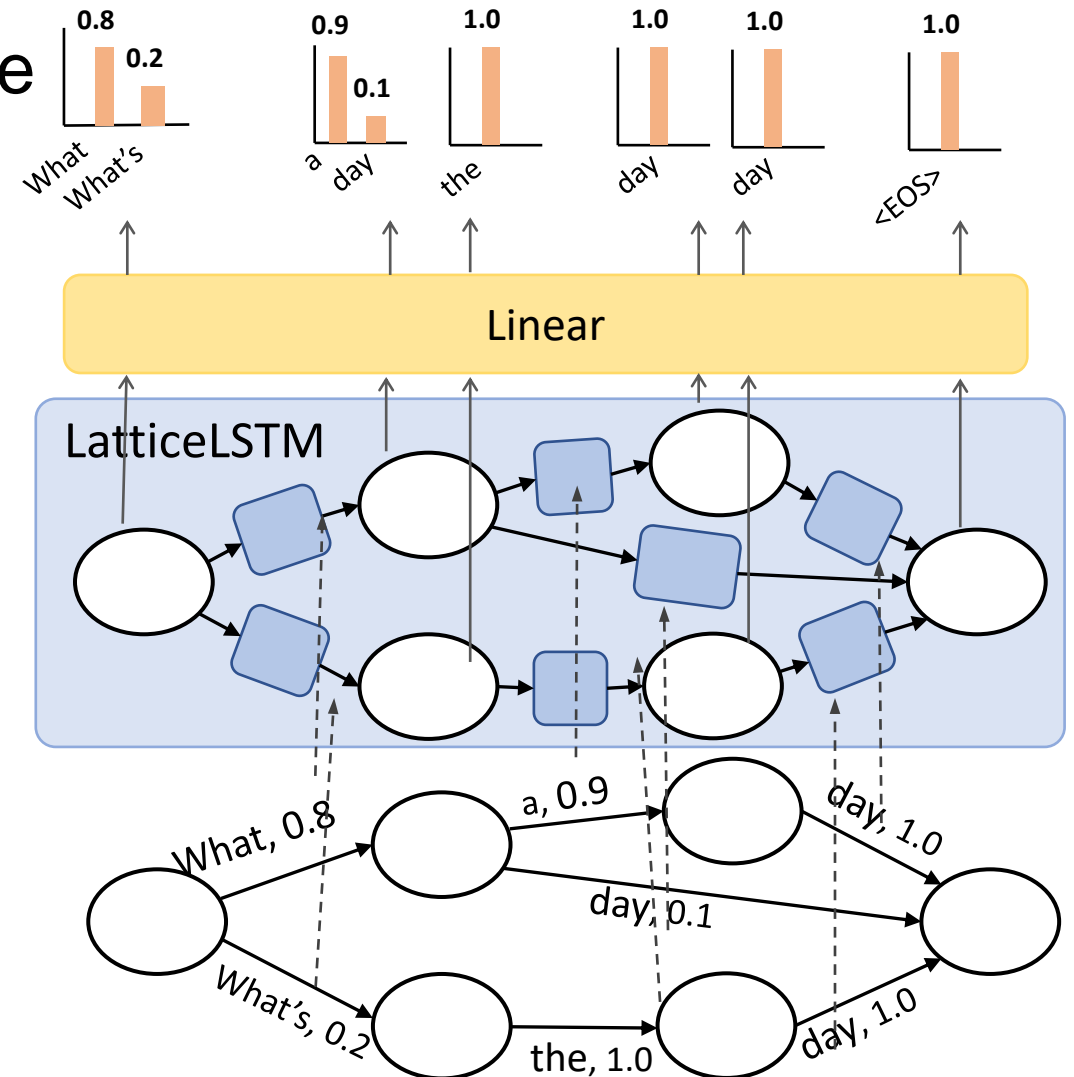
- Idea: lattices may include correct words



Lattice Language Modeling

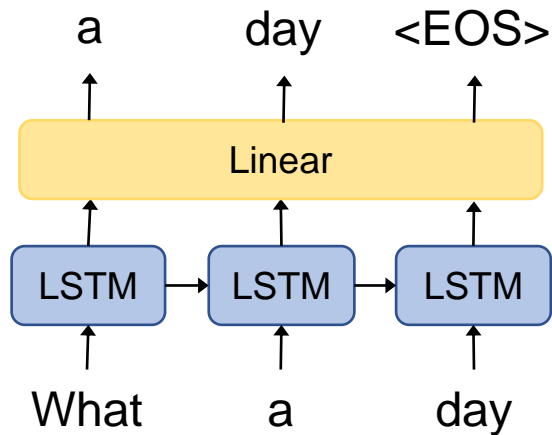
- 1) LatticeLSTM encodes nodes of a lattice
 - 2) The goal is to predict the outgoing transitions (words) given a node's representation
- The one-hypothesis lattice reduces to normal language modeling

Issue: LatticeLSTM runs prohibitively slow

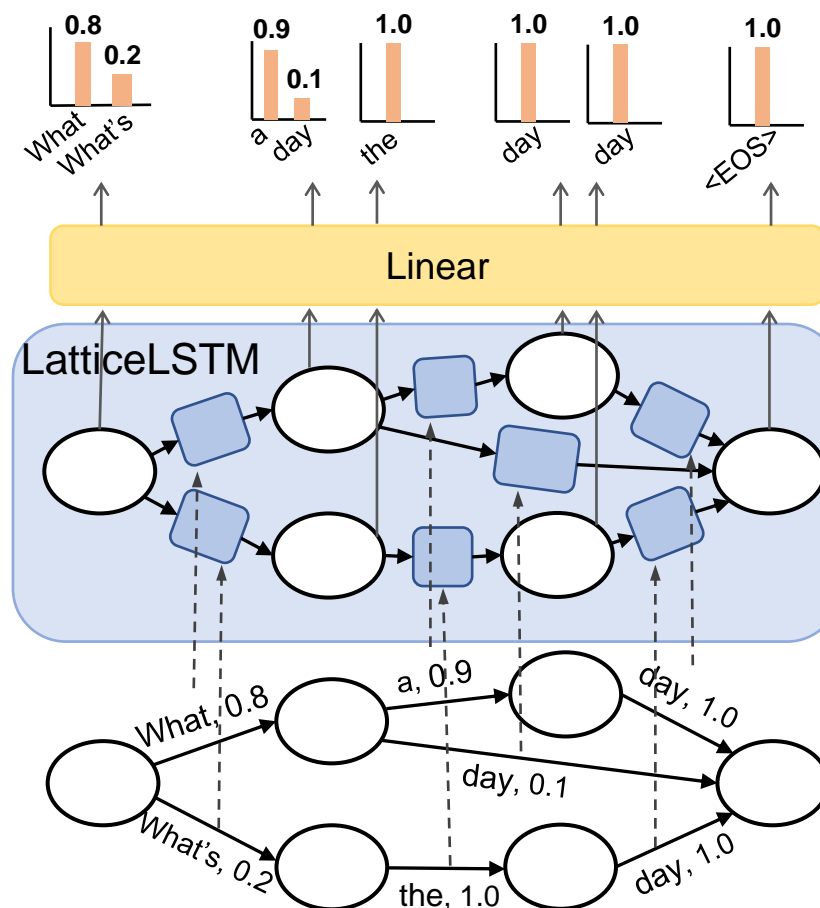


Efficient Two-Stage Pre-Training

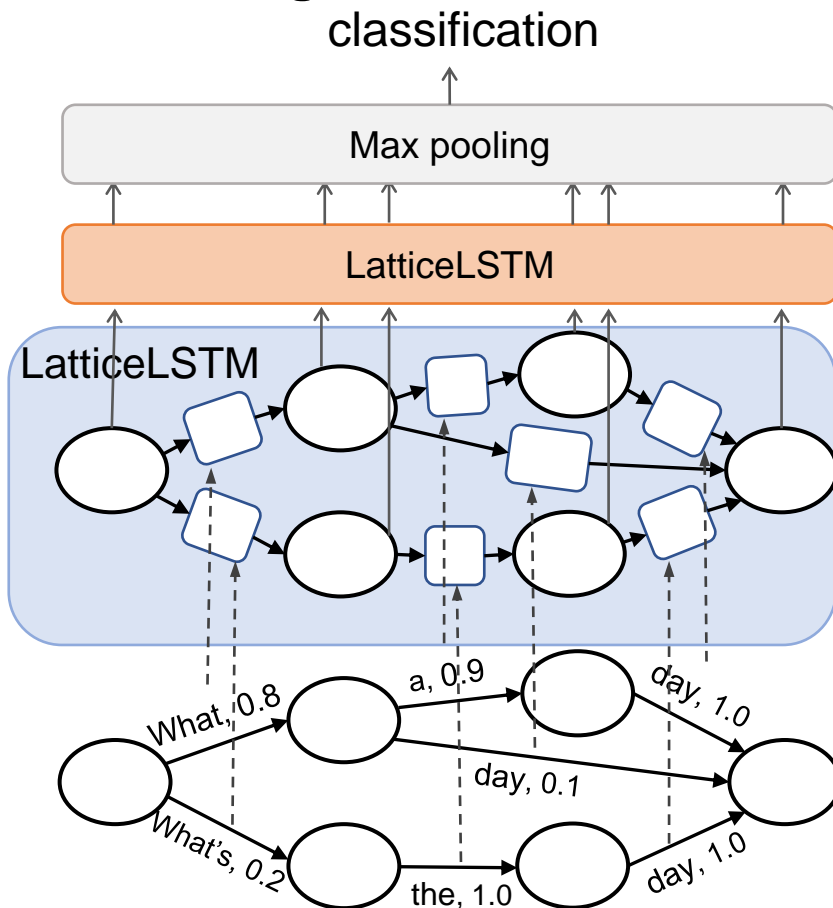
Stage 1: Pre-Training on Sequential Texts



Stage 2: Pre-Training on Lattices



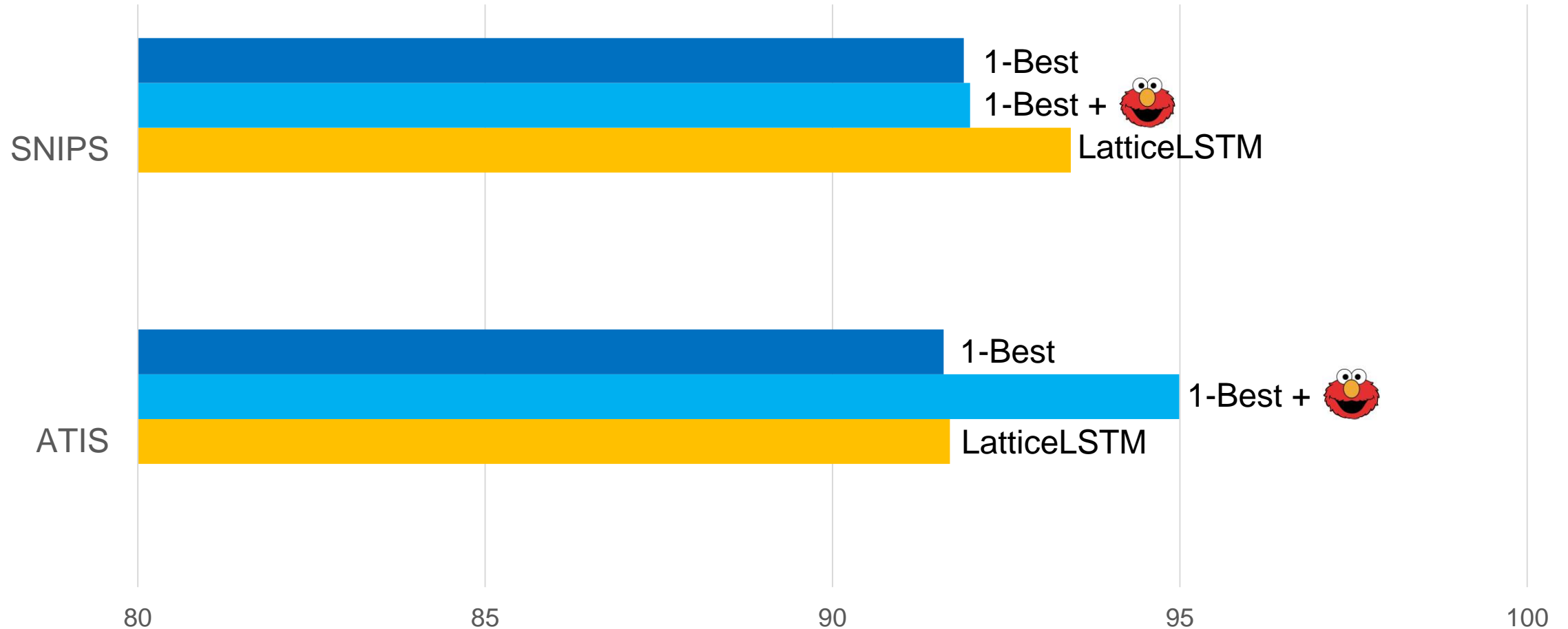
Fine-Tuning



Spoken Language Understanding Results

Intent Prediction

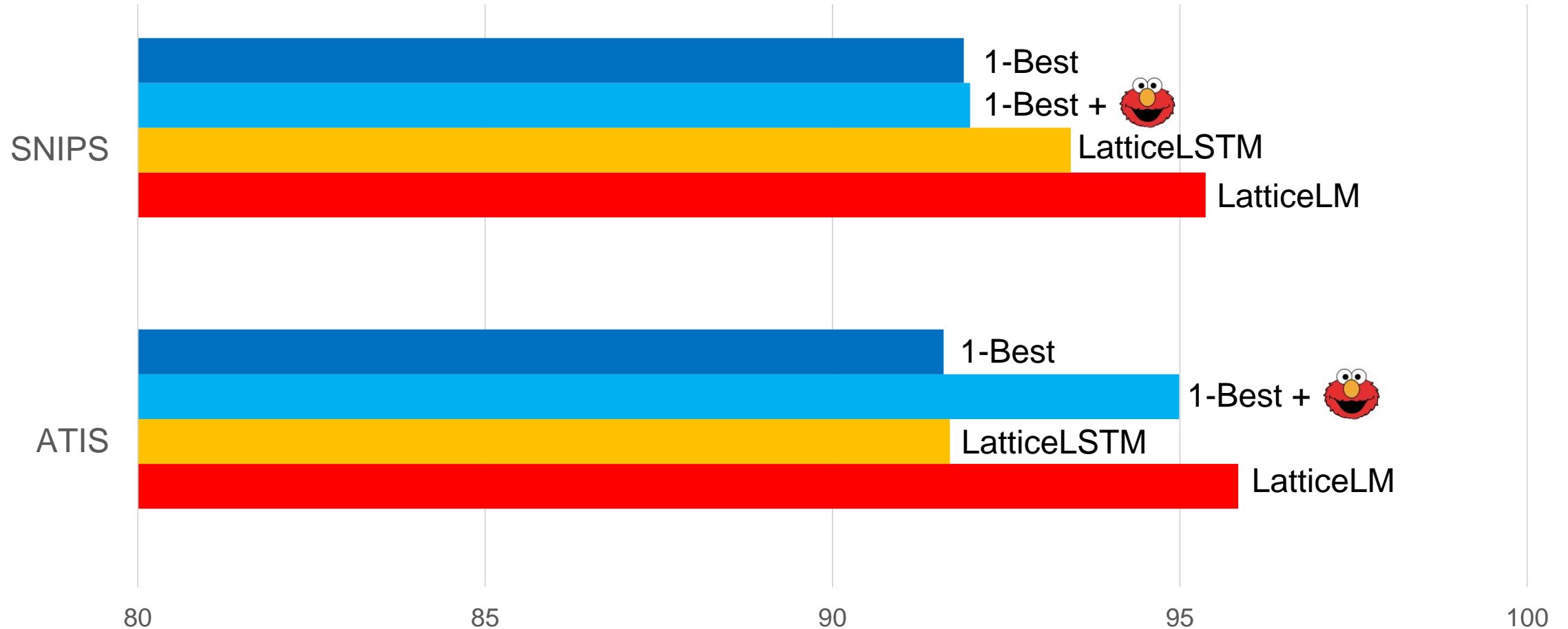
- Word Error Rate: 45.6% (SNIPS); 15.6% (ATIS)



Spoken Language Understanding Results

Intent Prediction

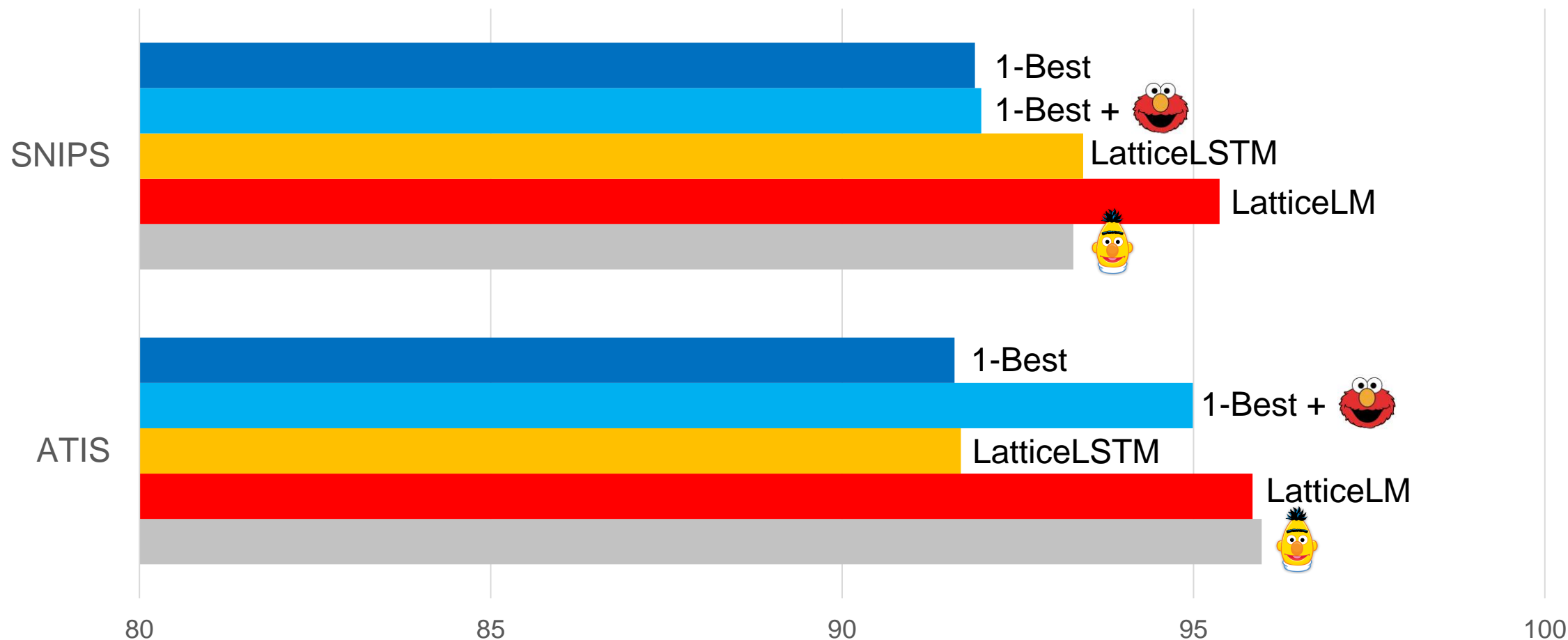
- Word Error Rate: 45.6% (SNIPS); 15.6% (ATIS)



Spoken Language Understanding Results

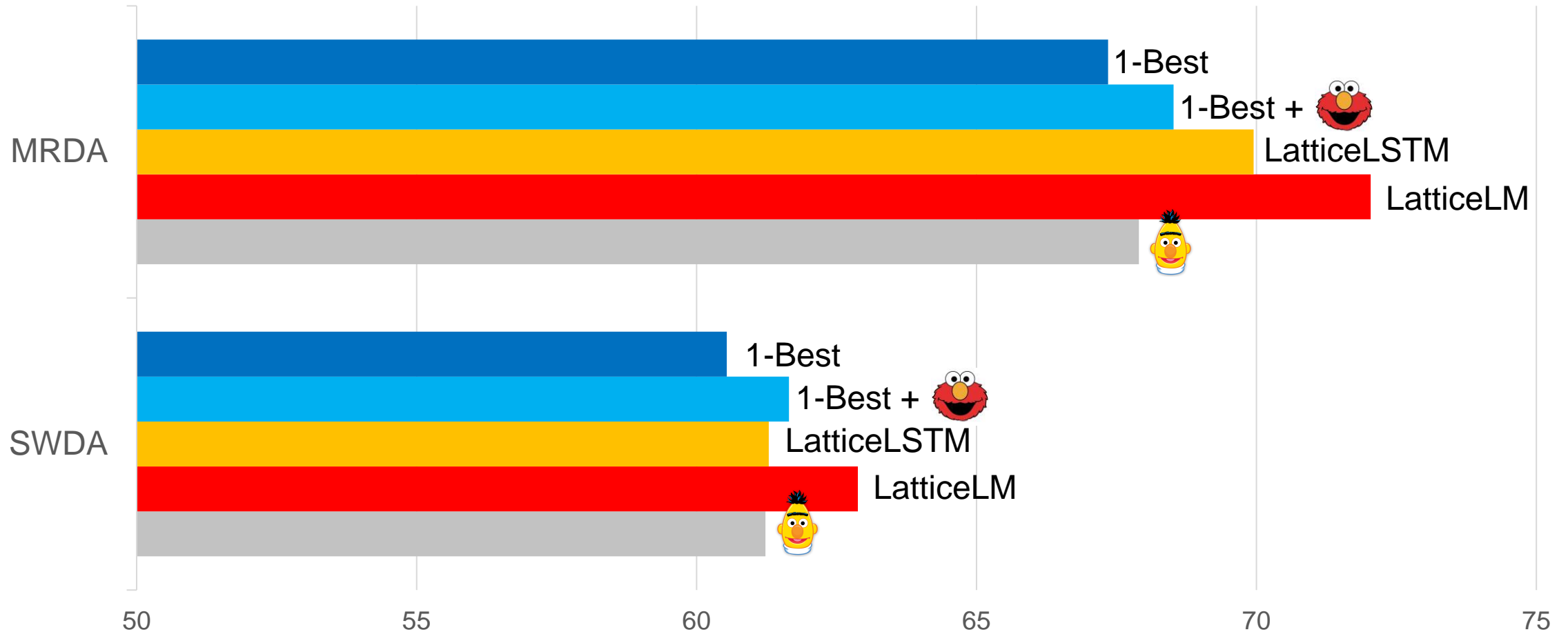
Intent Prediction

- Word Error Rate: 45.6% (SNIPS); 15.6% (ATIS)



Spoken Language Understanding Results

- Dialogue Act Prediction
 - Word Error Rate: 32.0% (MRDA); 28.4% (SWDA)



What if we do not have ASR lattices?

Solution:

Learning ASR-Robust Embeddings

(Huang & Chen, ICASSP 2020)

18

ASR-Robust Contextualized Embeddings

Confusion-Aware Fine-Tuning

Supervised

$$\text{Acoustic Confusion } \mathcal{C} = \{w_3^{x_{\text{trs}}}, w_2^{x_{\text{asr}}}\}$$

x_{trs} : Show me the fares from Dallas to Boston

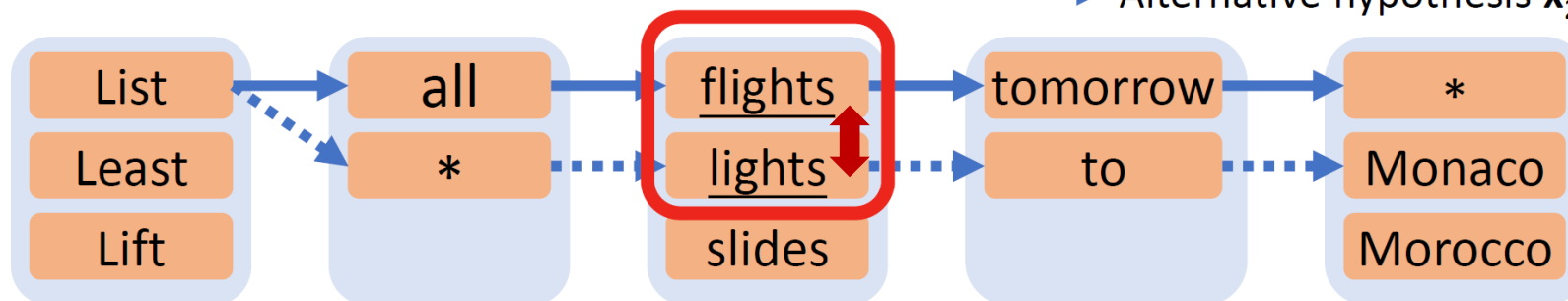
x_{asr} : Show me * affairs from Dallas to Boston



Unsupervised

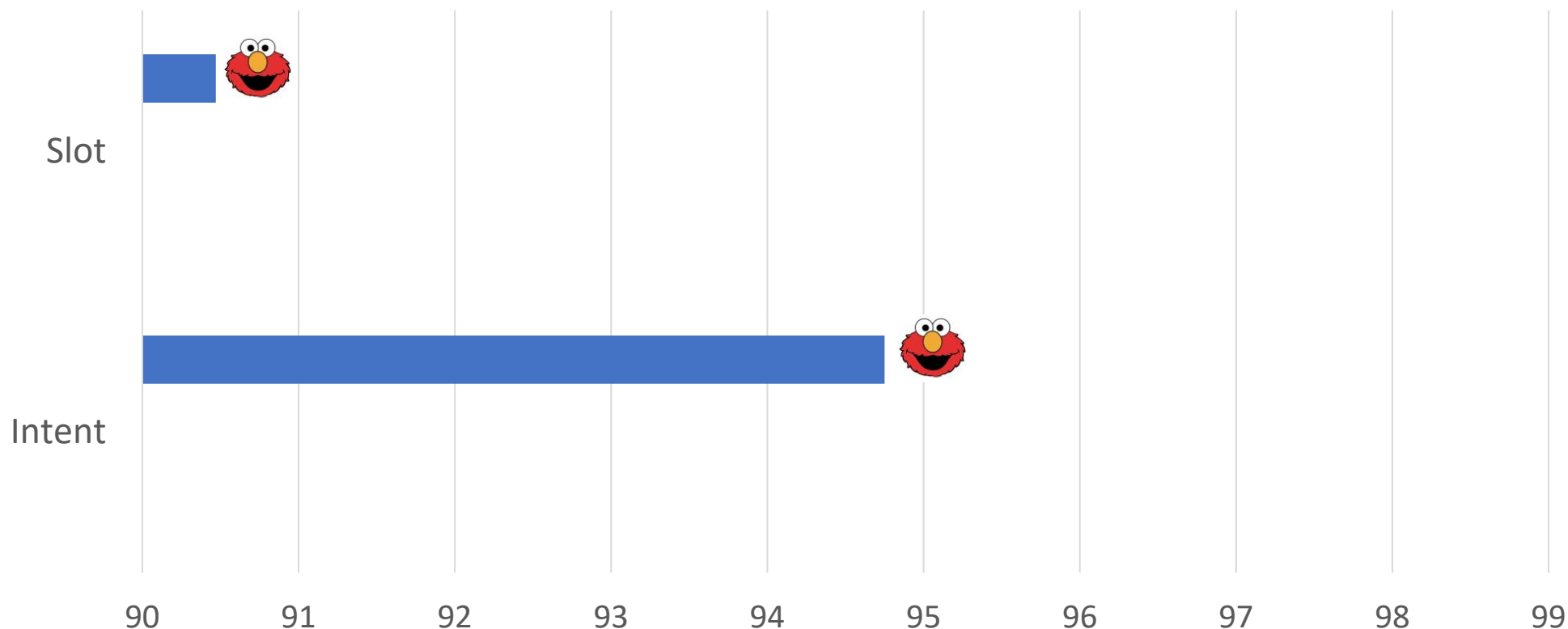
Acoustic Confusion

→ Top hypothesis x_1
 Alternative hypothesis x_2



Spoken Language Understanding Results

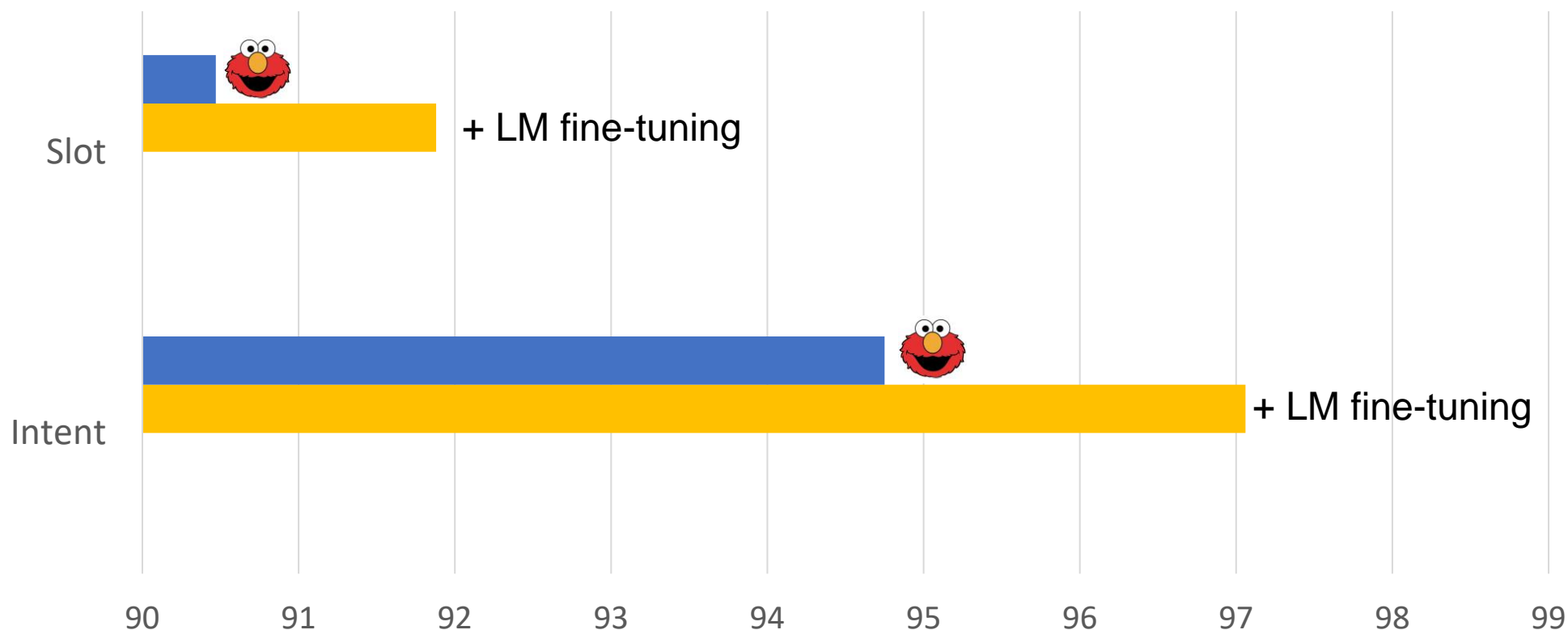
- Airline Traveling Information System (ATIS)
 - Word Error Rate: 16.4%



Chao-Wei Huang and Yun-Nung Chen, "Learning ASR-Robust Contextualized Embeddings for Spoken Language Understanding," in *The 45th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020.

Spoken Language Understanding Results

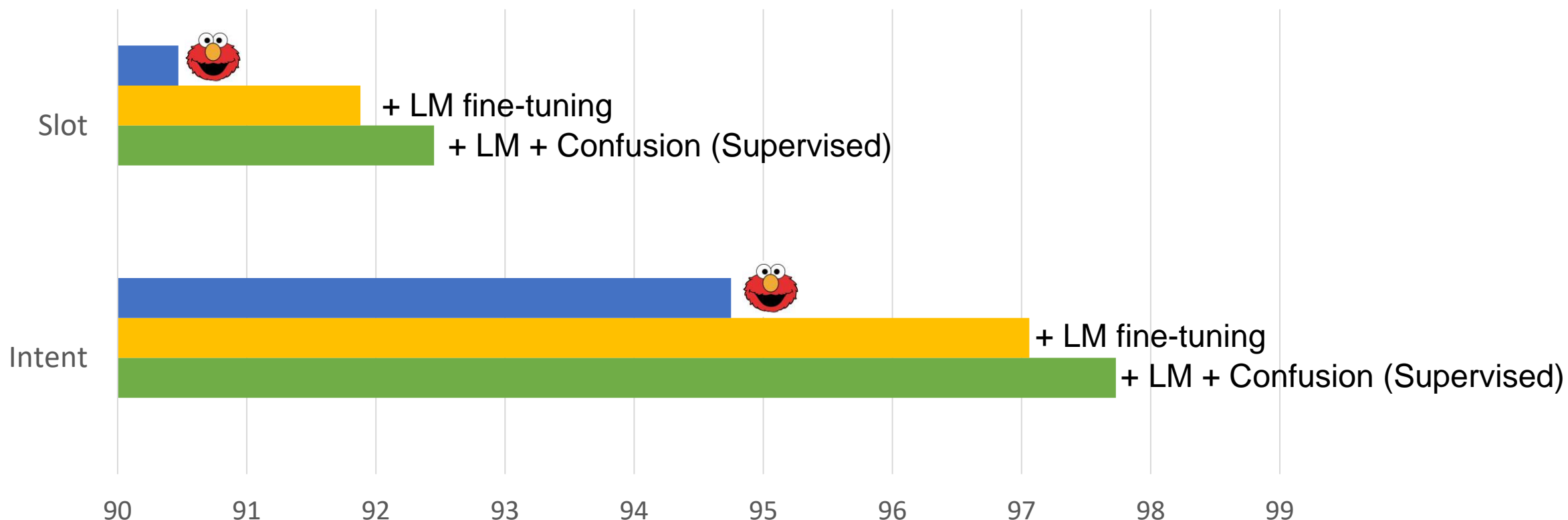
- Airline Traveling Information System (ATIS)
 - Word Error Rate: 16.4%



Chao-Wei Huang and Yun-Nung Chen, "Learning ASR-Robust Contextualized Embeddings for Spoken Language Understanding," in *The 45th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020.

Spoken Language Understanding Results

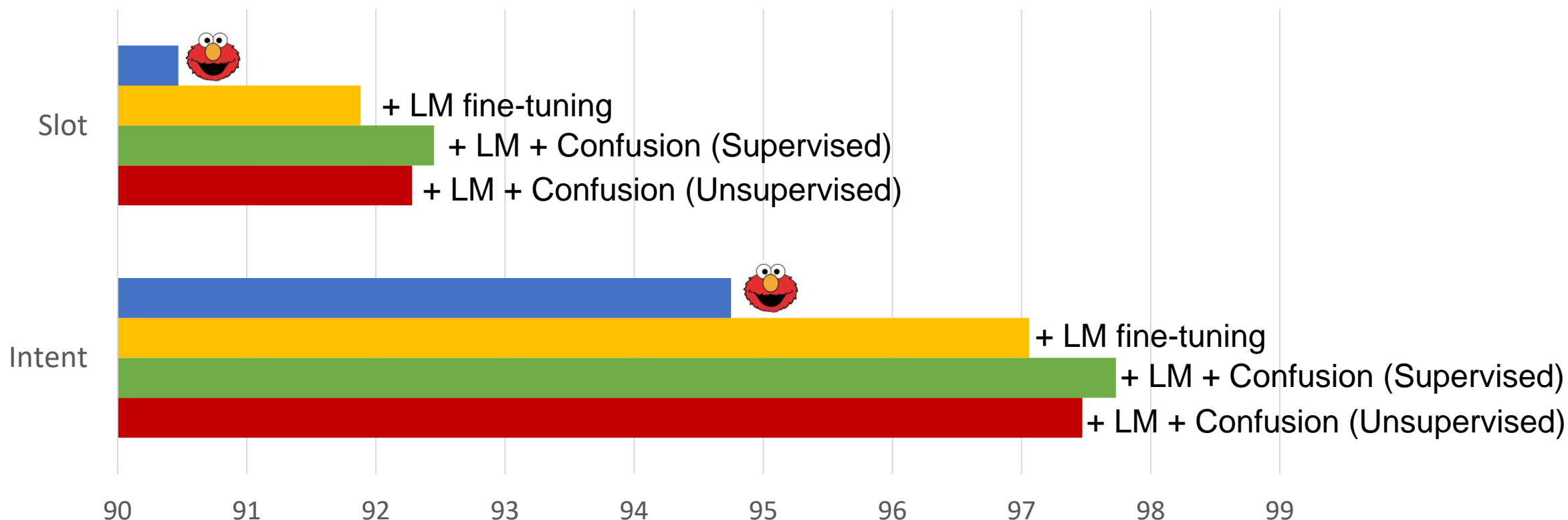
- Airline Traveling Information System (ATIS)
 - Word Error Rate: 16.4%



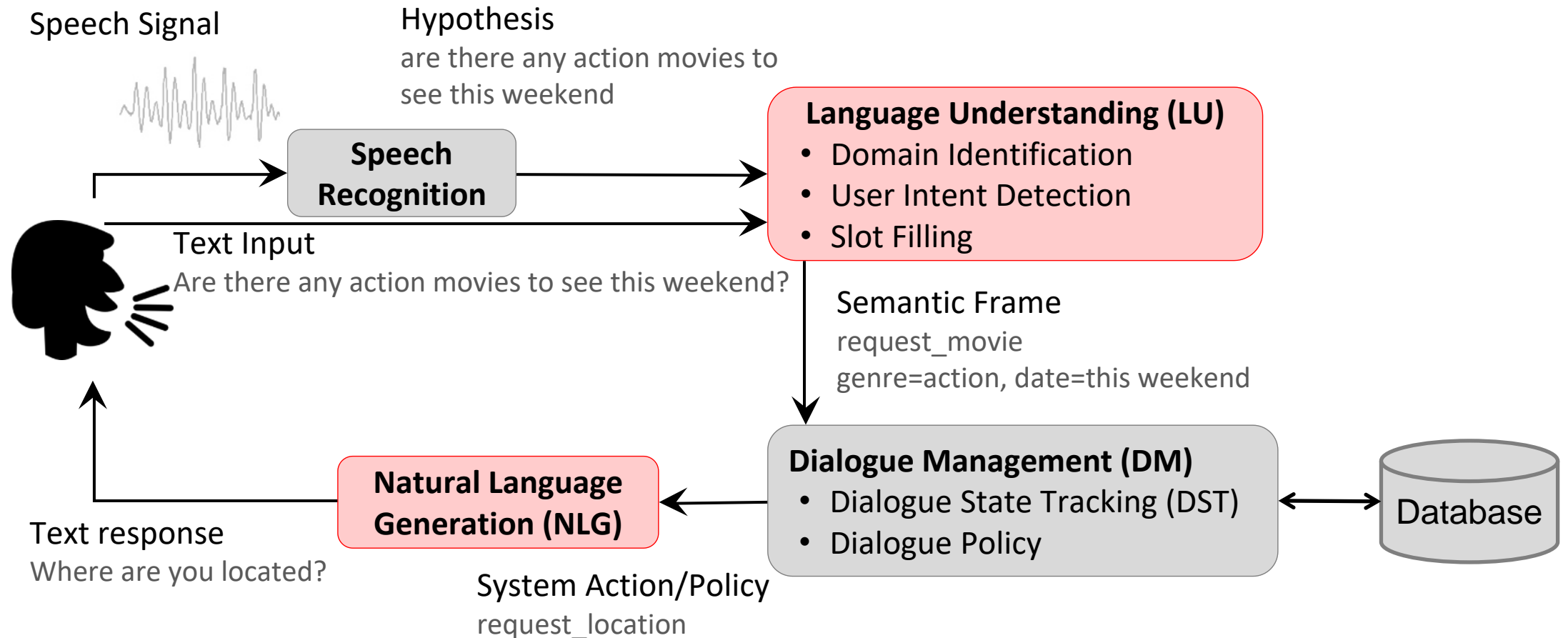
Chao-Wei Huang and Yun-Nung Chen, "Learning ASR-Robust Contextualized Embeddings for Spoken Language Understanding," in *The 45th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2020.

Spoken Language Understanding Results

- Airline Traveling Information System (ATIS)
 - Word Error Rate: 16.4%

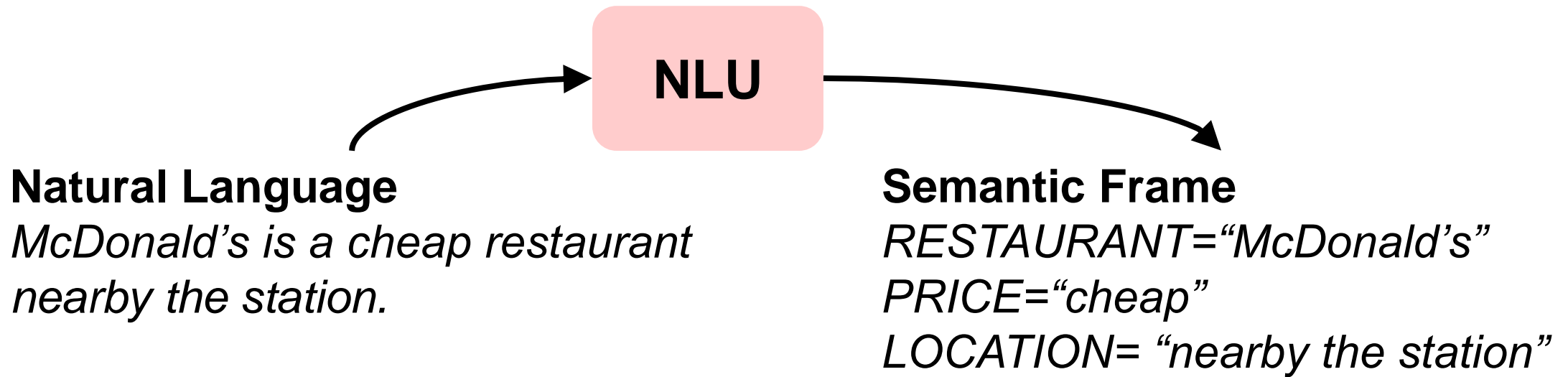


Task-Oriented Dialogue Systems ([Young, 2000](#))



Natural Language Understanding (NLU)

- Parse natural language into structured semantics



Natural Language Generation (NLG)

- Construct **natural language** based on **structured semantics**

Natural Language

*McDonald's is a cheap restaurant
nearby the station.*

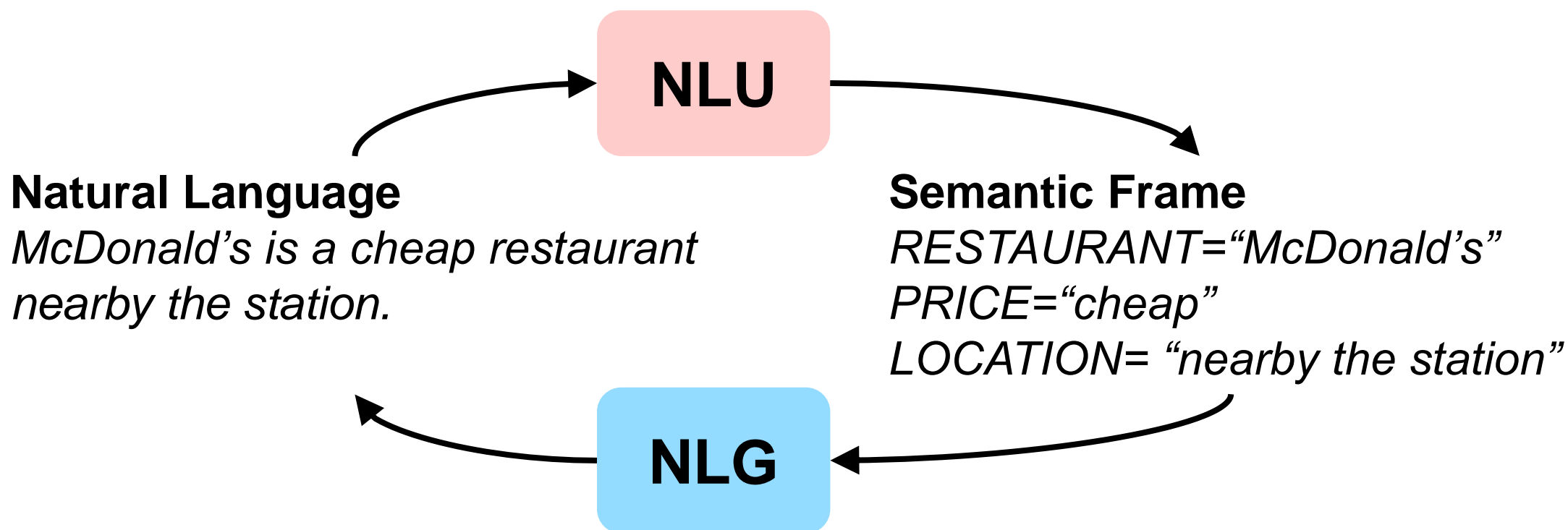
Semantic Frame

*RESTAURANT="McDonald's"
PRICE="cheap"
LOCATION="nearby the station"*

```
graph LR; SF[Semantic Frame] --> NLG[NLG]; NLG --> NL[Natural Language]
```

NLG

Duality between NLU and NLG



How can we leverage this dual relationship?

Solution:

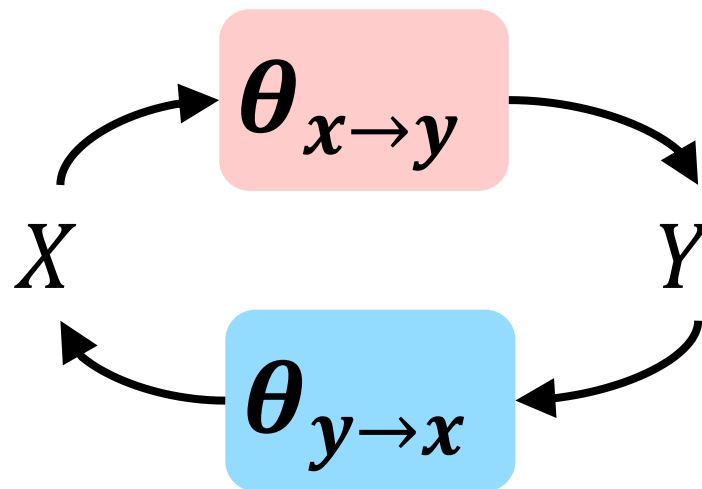
Dual Supervised Learning for NLU & NLG

(Su et al., ACL 2019)

28

DSL: Dual Supervised Learning (Xia et al., 2017)

- Proposed for machine translation
- Consider two domains X and Y , and two tasks $X \rightarrow Y$ and $Y \rightarrow X$



We have $P(x, y) = P(x | y)P(y) = P(y | x)P(x)$

Ideally $P(x, y) = P(x | y; \theta_{y \rightarrow x})P(y) = P(y | x; \theta_{x \rightarrow y})P(x)$

Dual Supervised Learning

- Exploit the duality by forcing models to follow the probabilistic constraint $P(x | y; \theta_{y \rightarrow x})P(y) = P(y | x; \theta_{x \rightarrow y})P(x)$

Objective function

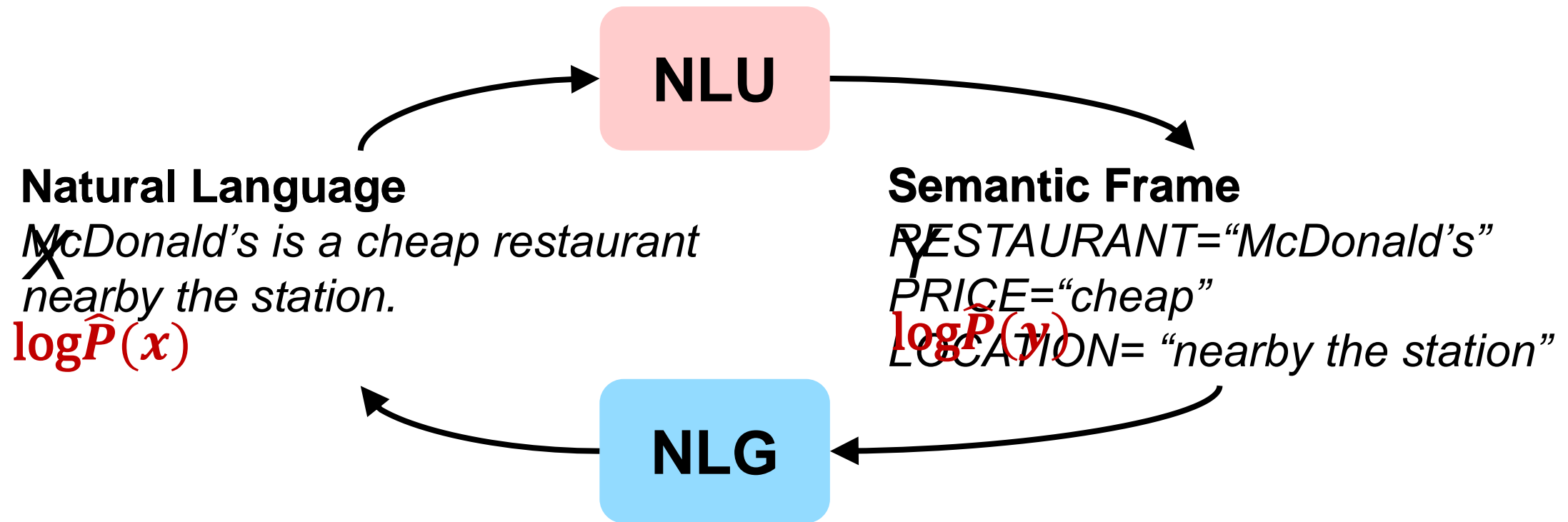
$$\begin{cases} \min_{\theta_{x \rightarrow y}} \mathbb{E}[l_1(f(x; \theta_{x \rightarrow y}), y)] + \lambda_{x \rightarrow y} l_{duality} \\ \min_{\theta_{y \rightarrow x}} \mathbb{E}[l_2(g(y; \theta_{y \rightarrow x}), x)] + \lambda_{y \rightarrow x} l_{duality} \end{cases}$$

$$l_{duality} = (\log \hat{P}(x) + \log P(y | x; \theta_{x \rightarrow y}) - \log \hat{P}(y) - \log P(x | y; \theta_{y \rightarrow x}))^2$$

How to model the marginal distributions of X and Y ?

Dual Supervised Learning

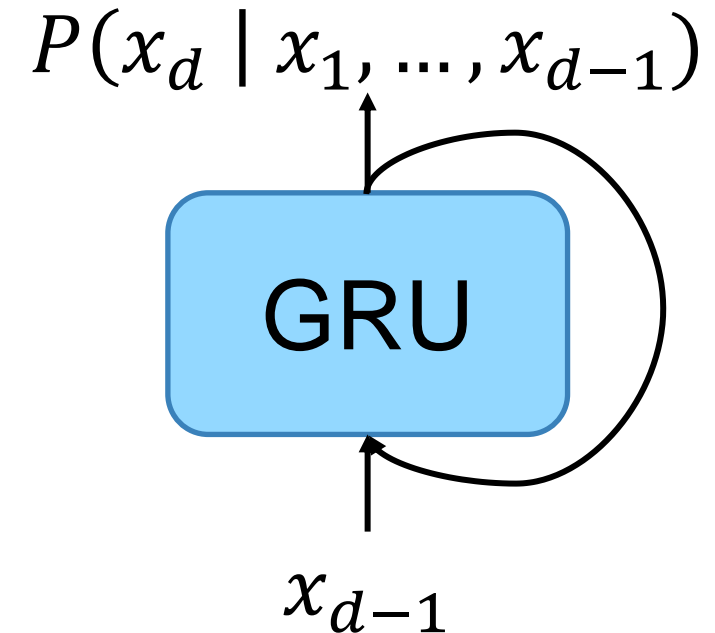
- Let's go back to NLU and NLG



Natural Language $\log \hat{P}(x)$

Language modeling

$$p(x) = \prod_d^D p(x_d \mid x_1, \dots, x_{d-1})$$



Semantic Frame $\log \hat{P}(y)$

- We treat NLU as a multi-label classification problem
- Each label is a slot-value pair

RESTAURANT="McDonald's"
PRICE="cheap"
LOCATION="nearby the station"



0
1
.
.
.
0
1

How to model the marginal distributions of y ?

Semantic Frame $\log \hat{P}(y)$

Naïve approach

- Calculate prior probability for each label $\hat{P}(y_i)$ on the training set.
- $\hat{P}(y) = \prod \hat{P}(y_i)$

Assumption: labels are independent

Restaurant: "McDonald's"	Price: "cheap"	Food: "Pizza"
Restaurant: "KFC"	Price: "expensive"	Food: "Hamburger"
Restaurant: "PizzaHut"		Food: "Chinese"

Semantic Frame $\log \hat{P}(y)$

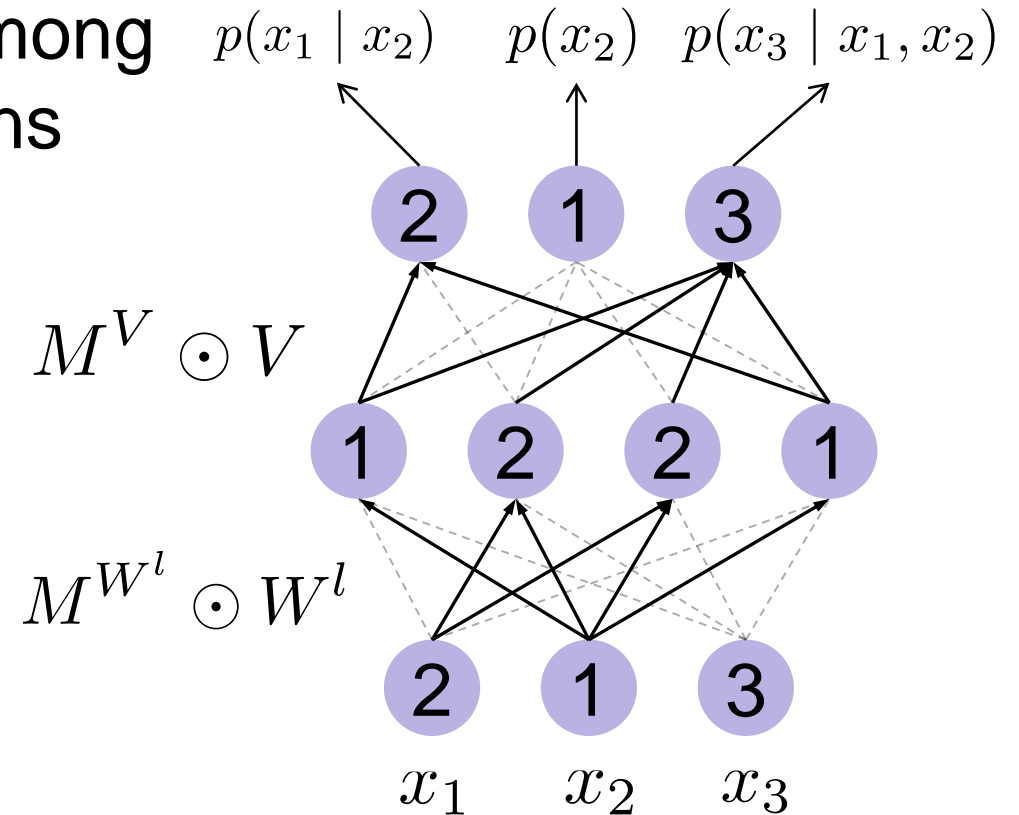
Masked autoencoder for distribution estimation (MADE)

Introduce sequential dependency among labels by masking certain connections

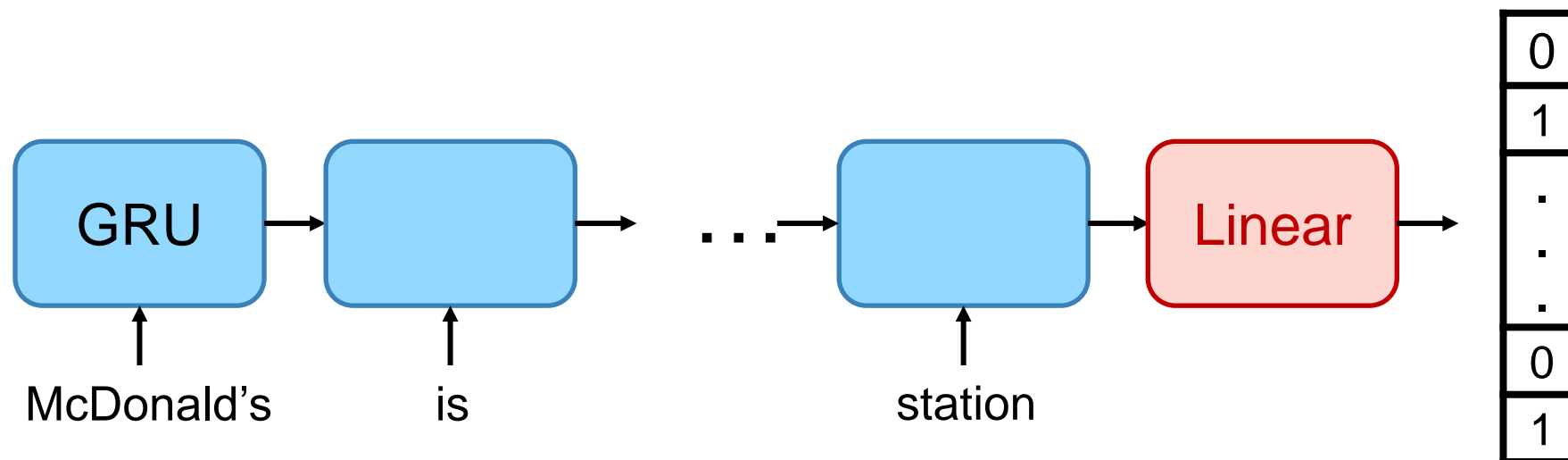
$$M = \begin{cases} 1 & \text{if } m^l(k') \geq m^{l-1}(k) \text{ or } m^L(d) > m^{L-1}(k) \\ 0 & \text{otherwise} \end{cases}$$

$$p(x) = \prod_d p(x_d \mid S_d)$$

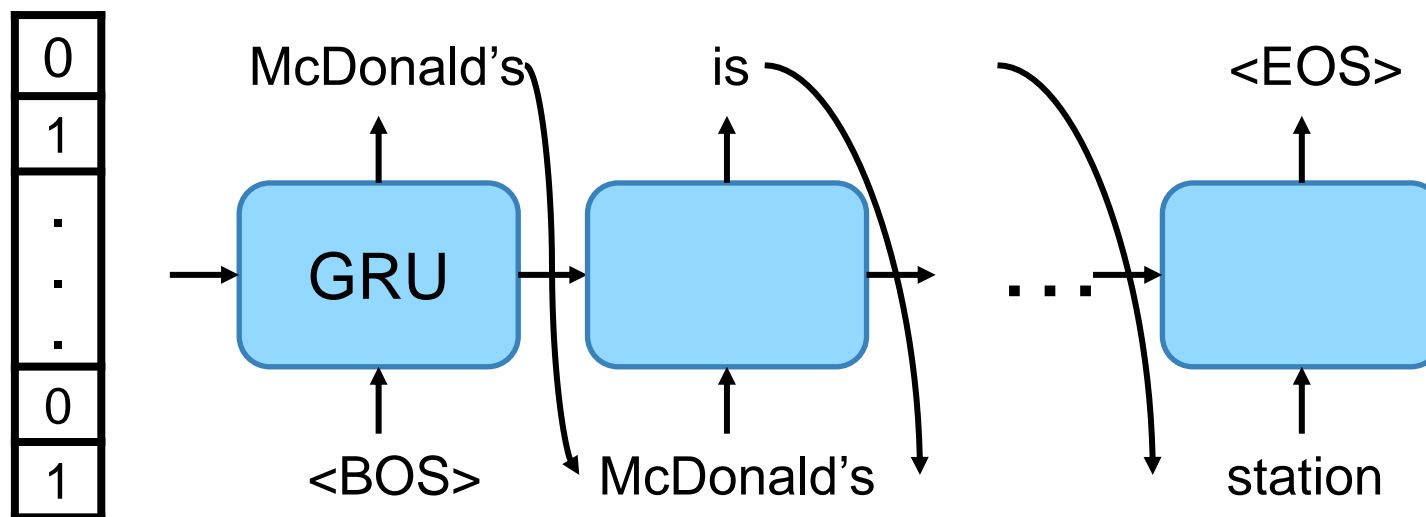
→ marginal distribution of y



NLU

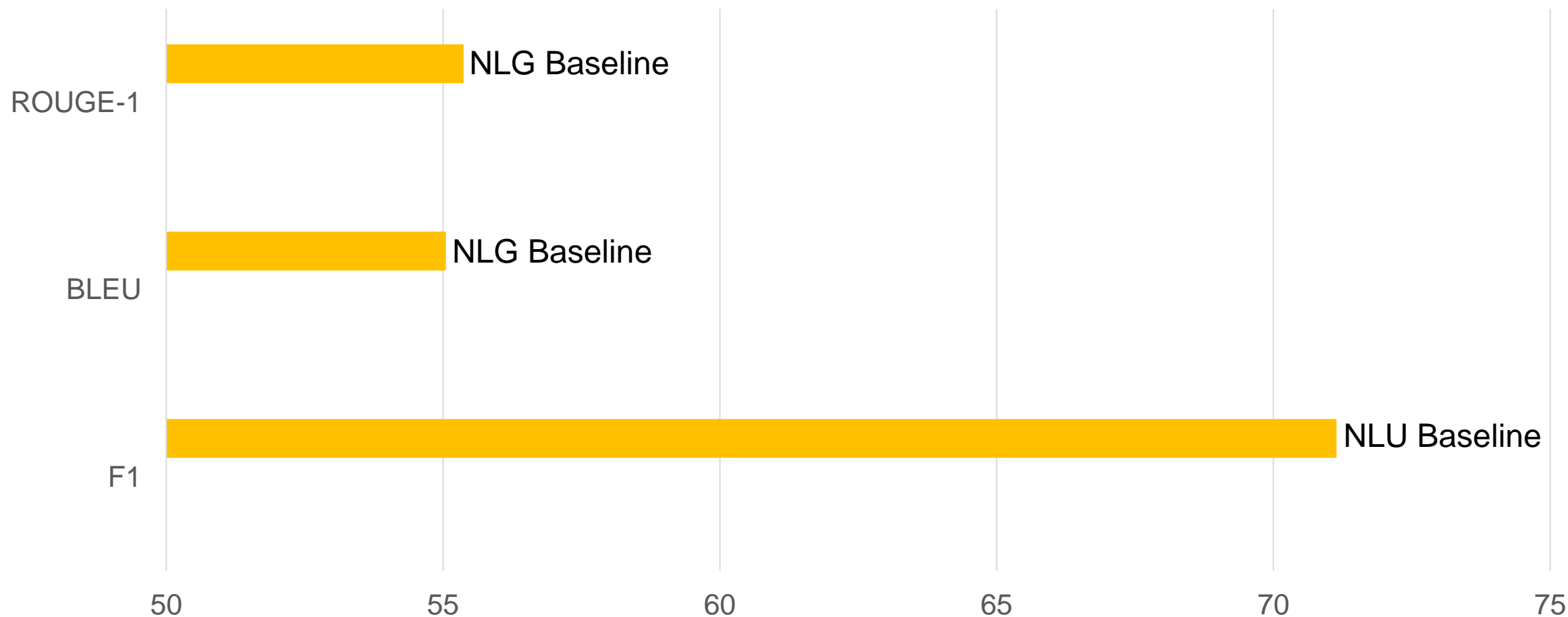


NLG



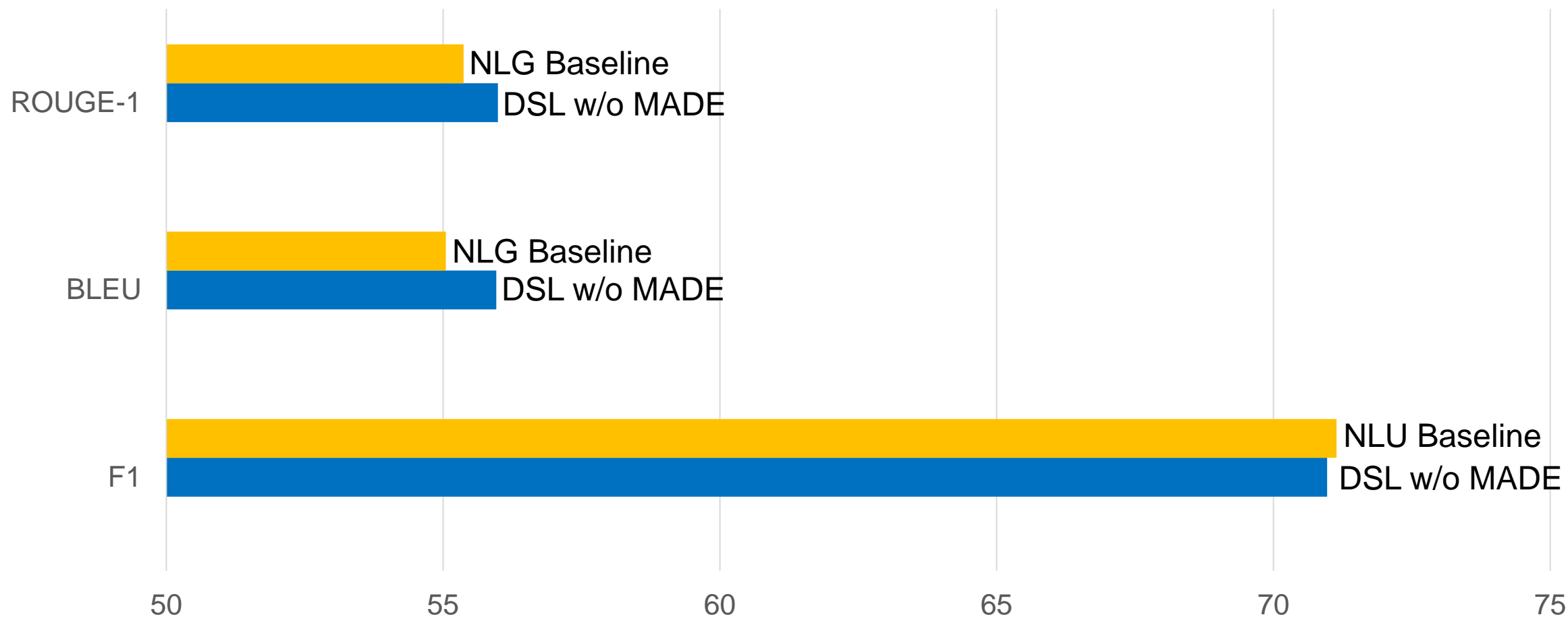
NLU/NLG Results

- E2E NLG data: 50k examples in the restaurant domain
- NLU: F-1 score; NLG: BLEU, ROUGE



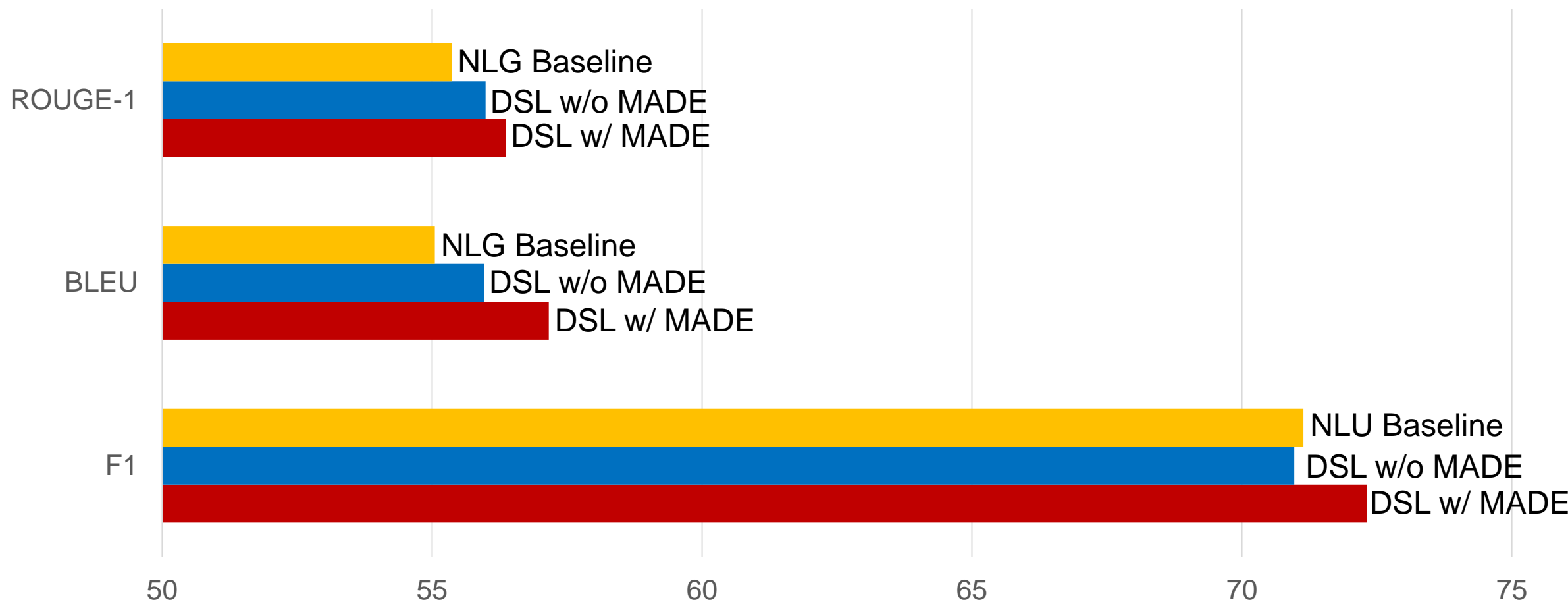
NLU/NLG Results

- E2E NLG data: 50k examples in the restaurant domain
- NLU: F-1 score; NLG: BLEU, ROUGE



NLU/NLG Results

- E2E NLG data: 50k examples in the restaurant domain
- NLU: F-1 score; NLG: BLEU, ROUGE



Summary

- **Robustness: spoken language embeddings** are needed for better conversational AI

- Written texts enough for pre-training embeddings
- Mismatch when applying to spoken language

1) LatticeLM for preserving uncertainty

2) Adapting contextualized embeddings robust to misrecognition



- **Scalability: leveraging the duality** of NLU and NLG

- Apply dual learning to leverage the duality
- Data distribution property is important
- Better performance and flexibility for diverse NLU/NLG models



- Yun-Nung (Vivian) Chen
- Assistant Professor, National Taiwan University
- y.v.chen@ieee.org / <http://vivianchen.idv.tw>

