#### 26 Task-Oriented Dialogue Systems (Young, 2000)



## 27—Natural Language Understanding (NLU)

Parse natural language into structured semantics



## 28—Natural Language Generation (NLG)

Construct natural language based on structured semantics

# Natural Language Semantic Frame McDonald's is a cheap restaurant RESTAURANT="McDonald's" pRICE="cheap" LOCATION= "nearby the station" NLG NLG

#### 29 Duality between NLU and NLG



# 30 Dual Supervised Learning for NLU & NLG (Su et al., 2019)

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

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



We have P(x, y) = P(x | y)P(y) = P(y | x)P(x)Ideally  $P(x, y) = P(x | y; \theta_{y \to x})P(y) = P(y | x; \theta_{x \to y})P(x)$ 

Xia, Y., Qin, T., Chen, W., Bian, J., Yu, N., & Liu, T. Y., "Dual supervised learning," in *Proc. of ICML*, 2017.

## **32**—Dual Supervised Learning

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

**Objective function** 

$$\begin{cases} \min_{\theta_{x \to y}} \mathbb{E} [l_1(f(x; \theta_{x \to y}), y)] + \lambda_{x \to y} \ l_{duality} \\ \min_{\theta_{y \to x}} \mathbb{E} [l_2(g(y; \theta_{y \to x}), x)] + \lambda_{y \to x} \ l_{duality} \\ l_{duality} = \left( \log \hat{P}(x) + \log P(y \mid x; \theta_{x \to y}) - \log \hat{P}(y) - \log P(x \mid y; \theta_{y \to x}) \right)^2 \end{cases}$$

#### How to model the marginal distributions of *X* and *Y*?

Xia, Y., Qin, T., Chen, W., Bian, J., Yu, N., & Liu, T. Y., "Dual supervised learning," in *Proc. of ICML*, 2017.

## 33 Dual Supervised Learning

Let's go back to NLU and NLG



## **34**—Natural Language $\log \hat{P}(x)$

#### Language modeling

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





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



#### 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 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

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

$$p(x) = \prod_{d}^{D} p(x_d \mid S_d)$$

 $\rightarrow$  marginal distribution of y



Germain, M., Gregor, K., Murray, I., & Larochelle, H., "MADE: Masked autoencoder for distribution estimation," in *Proceedings of International Conference on Machine Learning*, 2015.





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





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





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



## 42—Task-Oriented Dialogue Systems (Young, 2000)







## 45 Unstructured Knowledge Access

#### • A machine reads big text data

- serves as a teacher
- A user can ask questions
  - serves as a student
  - in a conversational manner

#### $\rightarrow$ Conversational QA

#### Section: Sectio: Section: Section: Section: Section: Section: Section: Sect **STUDENT: What is the origin of Daffy Duck?** TEACHER: $\hookrightarrow$ first appeared in Porky's Duck Hunt STUDENT: What was he like in that episode? TEACHER: $\hookrightarrow$ assertive, unrestrained, combative STUDENT: Was he the star? TEACHER: $\rightarrow$ No, barely more than an unnamed bit player in this short STUDENT: Who was the star? TEACHER: $\checkmark$ No answer STUDENT: Did he change a lot from that first episode in future episodes? TEACHER: $\hookrightarrow$ Yes, the only aspects of the character that have remained consistent (...) are his voice characterization by Mel Blanc STUDENT: How has he changed? TEACHER: $\hookrightarrow$ Daffy was less anthropomorphic STUDENT: In what other ways did he change? TEACHER: $\hookrightarrow$ Daffy's slobbery, exaggerated lisp (...) is barely noticeable in the early cartoons. STUDENT: Why did they add the lisp? TEACHER: $\hookrightarrow$ One often-repeated "official" story is that it was modeled after producer Leon Schlesinger's tendency to lisp. STUDENT: Is there an "unofficial" story? TEACHER: $\hookrightarrow$ Yes, Mel Blanc (...) contradicts that conventional belief . . .



## Solution: FlowDelta (Yeh & Chen, 2019)

## FlowDelta: Information Gain in Dialogue Flow

Idea: model the difference of hidden states in multi-turn dialogues



**Conversation Flow (over Context)** 



#### 48 FlowDelta (Yeh & Chen, 2019)

Idea: model the difference of hidden states in multi-turn dialogues





#### Data: QuAC, CoQA





#### Data: QuAC, CoQA





#### Data: QuAC, CoQA



#### **52**—QuAC Leaderboard

Rank	Model	F1	HEQQ	HEQD
	Human Performance (Choi et al. EMNLP '18)	81.1	100	100
Sep 15, 2019	History-Attentive-TransBERT (single model) Alibaba Al Labs	72.9	69.7	13.6
2 Aug 31, 2019	<b>TransBERT (single model)</b> Anonymous	71.4	68.1	10.0
3 Apr 24, 2019	Bert-FlowDelta (single model) National Taiwan University, MiuLab https://arxiv.org/abs/1908.05117	67.8	63.6	12.1
4 June 13, 2019	Context-Aware-BERT (single model) Anonymous	69.6	65.7	8.1
5 Aug 22, 2019	BertMT (single model) WeChat Al	68.9	65.2	8.9
6 Sep 9, 2019	BertInfoFlow (single model) PINGAN Omni-Sinitic	69.3	65.2	8.5



- Spoken language embeddings are needed for better conversational AI
  - Written texts enough for pre-training embeddings
  - Mismatch when applying to spoken language
- 1) Adapting Transformer to ASR lattices



- 2) Adapting contextualized embeddings robust to misrecognition
- Leveraging the duality of NLU and NLG improves the scalability
  - Apply dual supervised learning to leverage the duality
  - Data distribution property is important
  - Better performance and flexibility for diverse NLU/NLG models

Conversational QA enables unstructured information access

• FlowDelta: information gain in dialogue flow guides better understanding