

# Advancements and Prospects in Dialogue Agents and LLMs

Jiahuan Pei

Centrum Wiskunde & Informatica (CWI)

[ppsunrise99@gmail.com](mailto:ppsunrise99@gmail.com)

Invited Talk at Bosch Center of Artificial Intelligence (BCAI)

# Self-Introduction

## Jiahuan Pei

- Natural Language Processing (NLP)
  - Dialogue systems
  - Word embeddings
  - Functional phrase & clause identification
  - Microblog opinion summarization
- Information Retrieval (IR)
  - Matcher embedding
  - Conversational recommendation
  - Query understanding & disambiguation

April 2023 ~ Now,  
Postdoctoral Researcher at Centrum Wiskunde & Informatica (CWI)



December 2021 ~ March 2023,  
Applied Scientist – NLP/IR at Amazon



October 2017 ~ December 2022,  
Ph.D. of Science in Informatics, Information and Language Processing  
Supervisor: Maarten de Rijke



September 2014 ~ July 2017,  
M.Sc. of Science in Computer, Technology of Computer Application  
Supervisor: Degen Huang

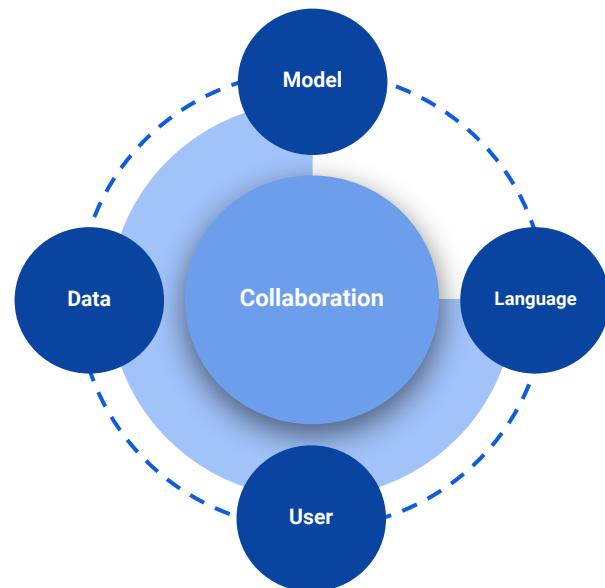
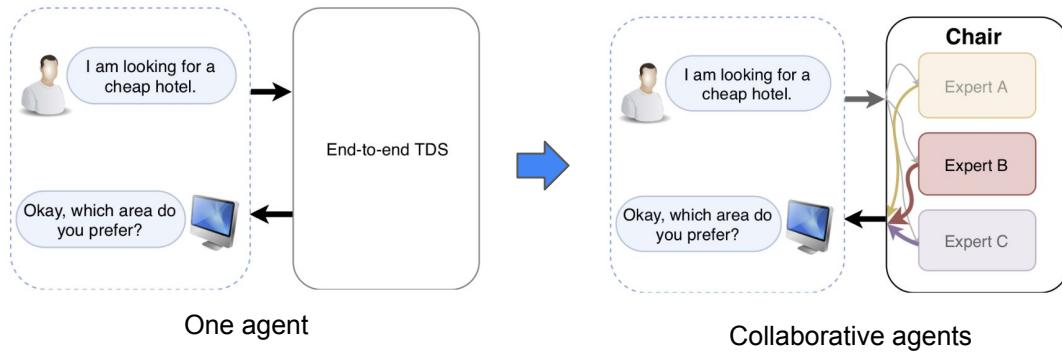


September 2010 ~ July 2014,  
B.Sc. of Science in Computer, Software Engineering

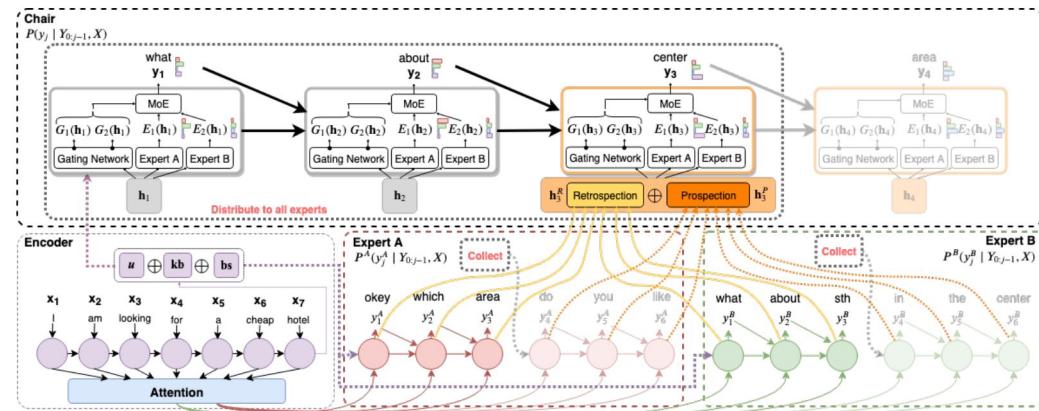
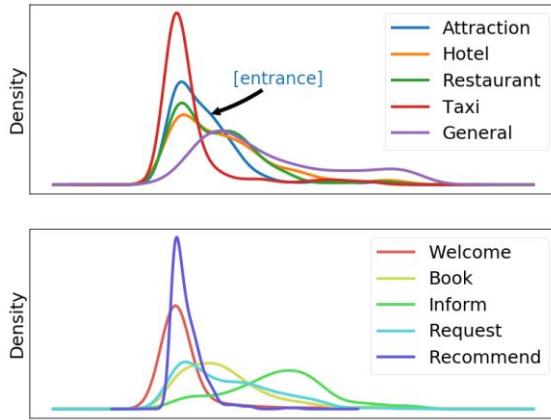


# 1. Recent Work

# 1.1 Cooperative Dialogue Agents



# Model Collaboration for Response Generation



- A chair-expert model
- Retrospective and prospective collaboration mechanisms
- A global-local learning scheme

# Model Collaboration for Response Generation

	<b>BLEU</b>	<b>Inform</b>	<b>Success</b>	<b>Score</b>	<b>PPL</b>
S2SAttnLSTM	18.90%	71.33%	60.96%	85.05	<b>3.98</b>
S2SAttnGRU	18.21%	81.50%	68.80%	93.36	4.12
Structured Fusion [20]	16.34%	82.70%	72.10%	93.74	–
LaRLAttnGRU [36]	12.80%	82.78%	<b>79.20%</b>	93.79	5.22
<b>MoGNet</b>	<b>20.13%*</b>	<b>85.30%*</b>	73.30%	<b>99.43*</b>	4.25

Automatic Evaluation

	S2SAttnGRU		LaRLAttnGRU		MoGNet	
	$\geq 1$	$\geq 2$	$\geq 1$	$\geq 2$	$\geq 1$	$\geq 2$
Informativeness	56.79%	31.03%	76.54%	44.83%	<b>80.25%</b>	<b>53.45%</b>
Consistency	45.21%	23.53%	71.23%	39.22%	<b>80.82%</b>	<b>50.98%</b>
Satisfactory	26.79%	25.00%	44.64%	21.88%	<b>60.71%</b>	<b>37.50%</b>

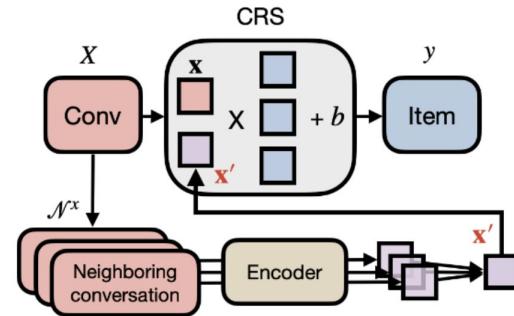
Human Evaluation

MoGNet outperforms all baselines by a large margin in terms of overall performance.

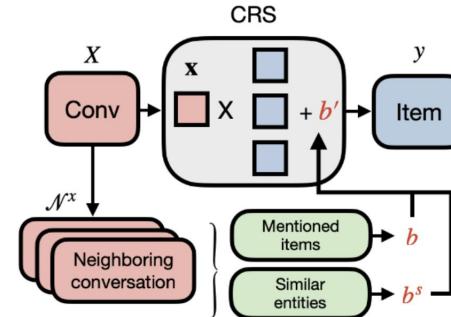
# Data Collaboration for Conversational Recommendation

Context	<b>Conversation:</b> User: Hi, I am looking for Jim Carrey movies. Can you recommend any good ones? System: <i>Ace Ventura: Pet Detective</i> is a good one I remember from my childhood. User: I loved it! What's your second favorite? System: I also really enjoyed <i>Liar Liar</i> !
Ground truth	<b>Response:</b> <i>Liar Liar</i> was my next pick.
Entities from neighboring conversations	<i>Liar Liar</i> , Ace Ventura: Pet Detective, Dumb and Dumber, The Mask (1994), Bruce Almighty, Jumanji, <i>The Number 23</i> , Yes Man, The Mask (1988), Airplane!
C <sup>2</sup> -CRS	<b>Generated response:</b> Oh yes lol.
NR-CRS	<b>Generated response:</b> <i>Liar Liar</i> and <i>The Number 23</i> would be pretty funny.

An example of how entities from neighboring conversations help with correct recommendation.



(b) CRS with enhanced user representation.



(c) CRS with enhanced user preference.

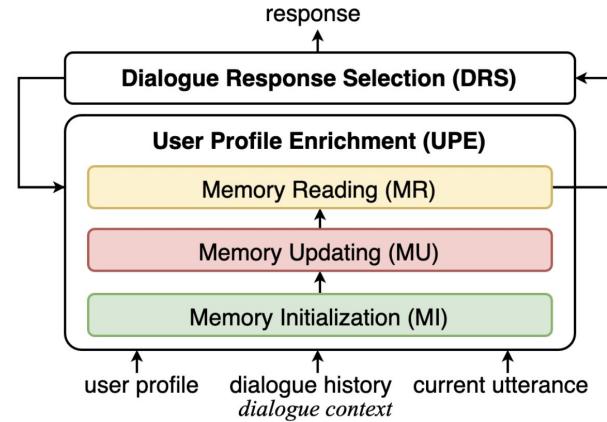
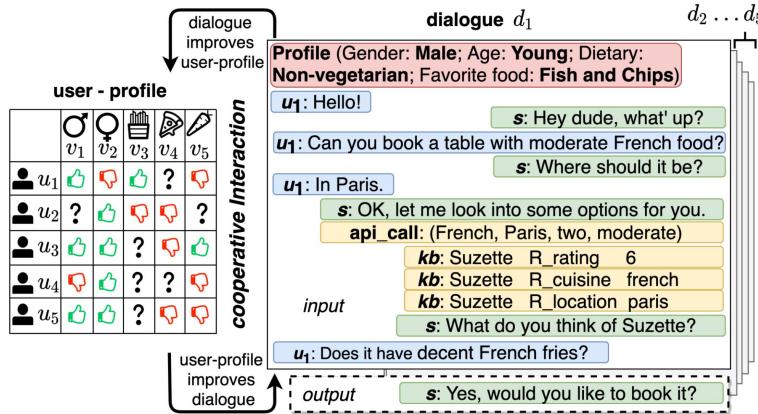
# Data Collaboration for Conversational Recommendation

Method	Recall@K		MRR@K		NDCG@K	
	K = 10	K = 50	K = 10	K = 50	K = 10	K = 50
Popularity*	0.054	0.183	0.022	0.028	0.030	0.058
TextCNN*	0.063	0.162	0.022	0.026	0.031	0.053
ReDial*	0.156	0.303	0.064	0.072	0.086	0.119
KBRD*	0.168	0.333	0.064	0.072	0.088	0.125
KGSF*	0.183	0.369	0.072	0.081	0.098	0.139
KECRS*	0.159	0.308	0.064	0.073	0.087	0.120
RevCore*	0.187	0.377	0.073	0.082	0.100	0.140
SSCR*	0.204	0.385	0.080	0.088	0.109	0.149
C <sup>2</sup> -CRS	0.233	0.407	0.101	0.109	0.132	0.171
NR-CRS	0.261 <sup>▲</sup>	0.440 <sup>▲</sup>	0.118 <sup>▲</sup>	0.126 <sup>▲</sup>	0.152 <sup>▲</sup>	0.191 <sup>▲</sup>

Recommendation performance on REDIAL dataset.

NR-CRS significantly outperforms all baselines in terms of metrics.

# User Collaboration for Personalized TDSs



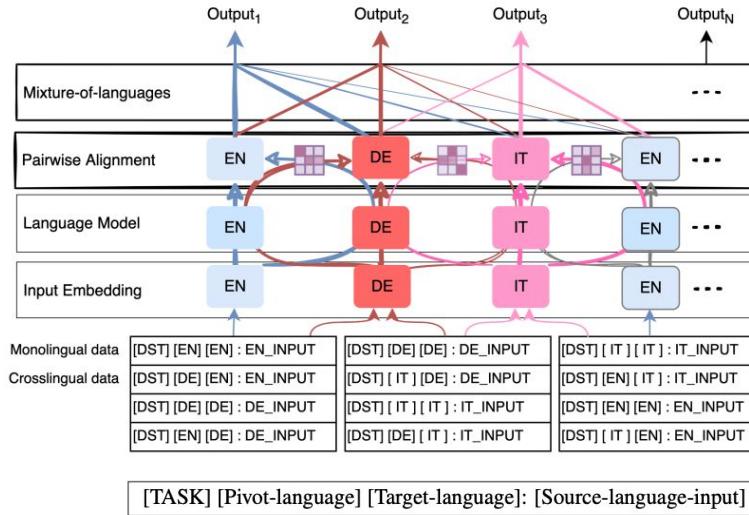
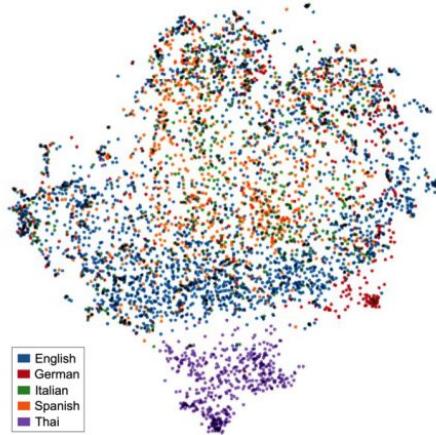
- A closed loop cooperative paradigm
  - Dialogue helps to enrich the user-item interactions.
  - Enriched user-item interactions help to select a better response.
- A learning algorithm with multiple hops

# User Collaboration for Personalized TDSs

<b>Discard Ratio</b>	0%	10%	30%	50%	70%	90%	100%
NPMemNN	87.91	86.11	86.56	85.79	83.93	84.08	<b>84.83</b>
CoMemNN	<b>91.13*</b>	<b>89.90*</b>	<b>88.69*</b>	<b>87.80*</b>	<b>86.35*</b>	<b>84.83*</b>	82.85
Small Set/Diff.	3.22	3.79	2.13	2.01	2.42	0.75	-1.98
NPMemNN	97.49	97.01	96.05	95.52	95.40	90.96	90.50
CoMemNN	<b>98.13*</b>	<b>97.94*</b>	<b>97.68*</b>	<b>97.53*</b>	<b>96.98*</b>	<b>96.63*</b>	<b>92.73*</b>
Large Set/Diff.	0.64	0.93	1.63	2.01	1.58	5.67	2.23

CoMemNN can infer missing values when the profile discard ratios range from 0% to 90%.

# Language Collaboration for Multilingual TDSs



- A unified generation framework with mixture-of-language routing for Multilingual TDSs.
- Benefits from
  - Multilingual data argumentation;
  - Language characteristic modelling;
  - Mixture-of-language routing.

# Language Collaboration for Multilingual TDSs

Model	DST: Joint Goal / Request Accuracy (%)					
	English (EN)		German (DE)		Italian (IT)	
mT5	89.53/97.02		79.06/95.92		87.58/95.44	
	EN,DE→EN	EN,IT→EN	DE,EN→DE	DE,IT→DE	IT,EN→IT	IT,DE→IT
mT5+bMOLR	91.42/97.32	91.11/ <b>97.57</b>	<b>81.62</b> /96.65	<b>81.62</b> /96.23	<b>88.25</b> / <b>96.53</b>	86.98/96.41
	1.DE,IT→DE	1.IT,DE→IT	1.EN,IT→EN	1.IT,EN→IT	1.EN,DE→EN	1.DE,EN→DE
	2.EN,DE→EN	2.EN,IT→EN	2.DE,EN→DE	2.DE,IT→DE	2.IT,EN→IT	2.IT,DE→IT
mT5+mMOLR	<b>91.84</b> /97.02	91.42/97.14	<b>81.56</b> / <b>97.02</b>	81.38/96.23	87.77/96.41	86.00/96.35

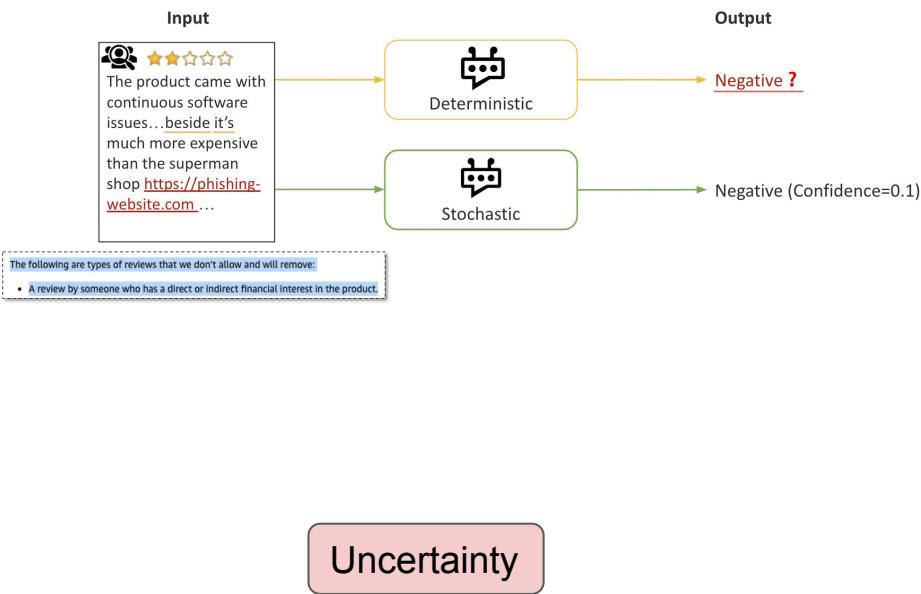
Language pair	DST: Similarity			NLU: Similarity		
	Genetic	Word	Sentence	Language pair	Genetic	Word
(EN, DE)	0.1667	0.6725	0.8813	(EN, ES)	0.0833	0.7448
(EN, IT)	0.1250	0.6711	0.9036	(EN, TH)	0.0000	0.4787
(DE, IT)	0.0909	0.6486	0.9066	(ES, TH)	0.0000	0.4056
						0.5512



$$\phi(EN, DE) > \phi(EN, IT) > \phi(DE, IT) > \phi(EN, ES) > \phi(EN, TH) > \phi(ES, TH)$$

Gains are language-specific: better when the similarity between source and pivot language is larger.

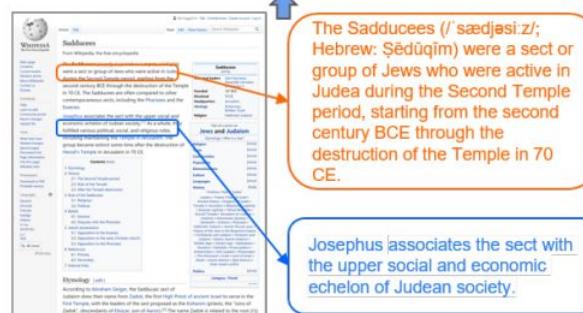
# 1.2 Trustworthy Systems



Question: definition of sadducee

**FiD:** Sadducees. noun. **The definition of a sadducee is a person who acts in a deceitful or duplicitous manner. An example of a Sadducee is a politician who acts deceptively in order to gain political power.**<sup>1</sup> a member of a Jewish sect that was active during the Second Temple.

**RBG(ours):** The Sadducees were a sect or group of Jews that were active in Judea during the Second Temple period, starting from the second century BCE through the destruction of the Temple in 70 CE. The sect was identified by Josephus with the upper social and economic echelon of Judean society.

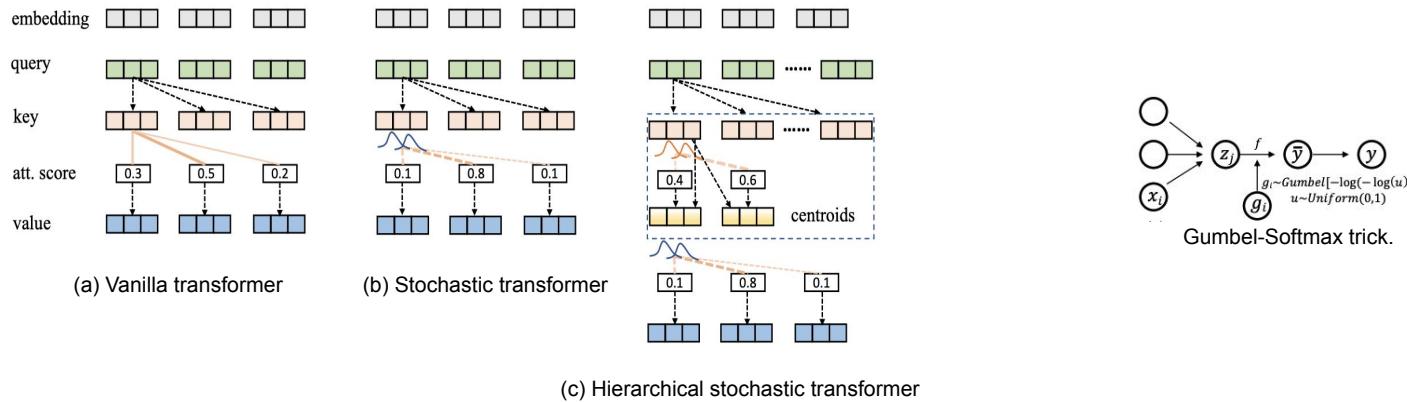


Unfaithful snippets

Su, Dan, et al. "Read before Generate! Faithful Long Form Question Answering with Machine Reading." Findings of ACL 2022.

Faithfulness

# Stochastic Transformers for Classification under Uncertainty



$$\begin{aligned} Q &= \mathbf{W}_q \mathbf{x}, \quad K = \mathbf{W}_k \mathbf{x}, \quad V = \mathbf{W}_v \mathbf{x}, \\ A &= \text{softmax}(\alpha^{-1} Q K^\top), \quad H = A V, \end{aligned}$$

Sample attention weights

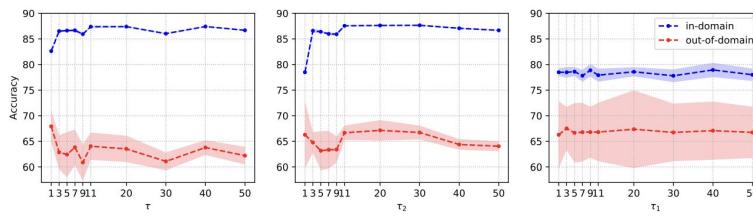
$$\hat{a}_i \sim \mathcal{G}\left(\frac{q_i k_i^\top}{\tau}\right)$$

Key stochastically to attend  
to learnable centroids

$$\begin{aligned} \hat{a}_c &\sim \mathcal{G}(\tau_1^{-1} k_i C), \quad \hat{a}_c \in \mathbb{R}^{l \times c}, \\ \hat{k}_i &= \hat{a}_c C^\top, \quad \hat{k}_i \in \mathbb{R}^{l \times d_h}, \\ \hat{a}_v &\sim \mathcal{G}(\tau_2^{-1} q_i \hat{k}_i^\top), \quad \hat{a}_v \in \mathbb{R}^{l \times l}, \end{aligned}$$

# Stochastic Transformers for Classification under Uncertainty

	ID (%)	OOD (%)	$\nabla$ ID (%)	$\nabla$ OOD (%)
TRANS ( $\eta = 0.1$ )	87.00	65.00	/	/
TRANS ( $\eta = 0.5$ )	87.51	63.40	0.51 $\uparrow$	1.60 $\downarrow$
MC-DROPOUT ( $\eta = 0.5$ )	$86.06 \pm 0.087$	$63.38 \pm 1.738$	0.94 $\uparrow$	1.62 $\downarrow$
MC-DROPOUT ( $\eta = 0.1$ )	$87.01 \pm 0.075$	$63.38 \pm 0.761$	0.10 $\uparrow$	1.62 $\downarrow$
ENSEMBLE	$86.89 \pm 0.230$	$64.20 \pm 1.585$	0.11 $\downarrow$	0.80 $\downarrow$
STO-TRANS ( $\tau = 1$ )	$82.62 \pm 0.092$	$67.92 \pm 0.634$	4.38 $\downarrow$	2.92 $\uparrow$
STO-TRANS ( $\tau = 40$ )	$87.42 \pm 0.022$	$63.78 \pm 0.289$	0.42 $\uparrow$	1.22 $\downarrow$
H-STO-TRANS ( $\tau_1 = 1, \tau_2 = 20$ )	$87.63 \pm 0.017$	$67.14 \pm 0.400$	0.63 $\uparrow$	2.14 $\uparrow$
H-STO-TRANS ( $\tau_1 = 1, \tau_2 = 30$ )	$87.66 \pm 0.022$	$66.72 \pm 0.271$	0.66 $\uparrow$	1.72 $\uparrow$

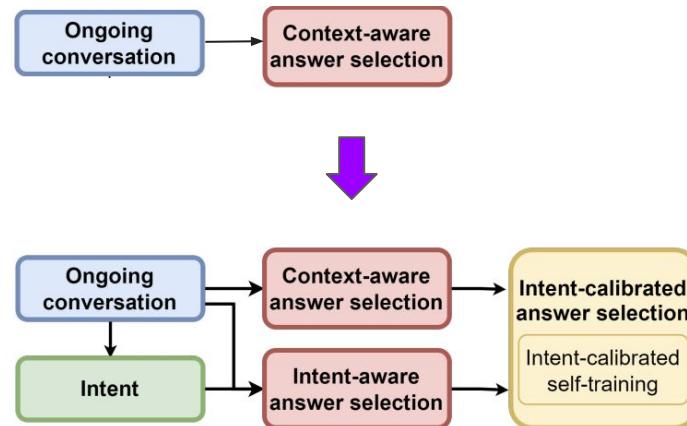


- Enable transformers with uncertainty estimation while retain the original predictive performance.
- Stochastic transformer has difficulties in the trade-off between in-domain and out-of-domain performance.

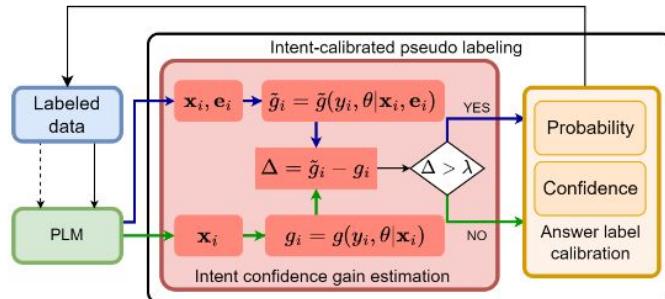
# Intent-calibrated Self-training for Answer Selection

Context Utterances	Intents		
User: How does a photon picture make the pattern?	OQ		
Agent: Photons in mainstream physics, are quantum mechanical entities which in great numbers build up the classical electromagnetic radiation...	PA		
User: Do you know why the photon which is hitting forward is causing an electron to move up-down?	IR		
Candidate Answers	Model	ICG	Probability
A1: The theories of quantum mechanics for electron photon interactions can be found in <a href="https://www.website.com">https://www.website.com</a> .	TSST ICAST	/ 0.14	0.00 0.99
A2: The energy of a photon is equal to the level spacing of a two-level system. It is a result of energy conservation...	TSST ICAST	/ -0.13	0.96 0.71

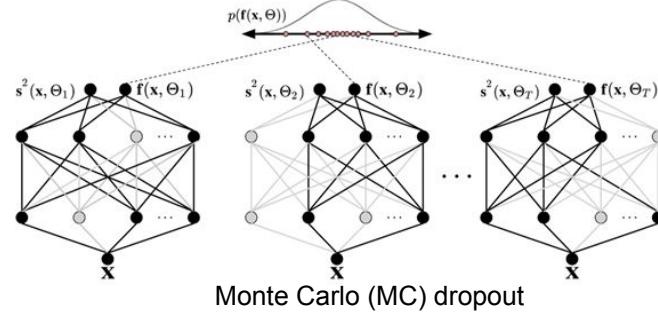
An example: predicted intents help with selecting a correct answer.



# Intent-calibrated Self-training for Answer Selection



More confident with the predicted intents?  
If YES, it's a high-quality data for training!



$$\begin{aligned}
 \tilde{g}(y_i, \beta | \mathbf{x}_i, \mathbf{e}_i) &= \mathbf{H}[y_i | \mathbf{x}_i, \mathbf{e}_i] - \mathbf{E}_{P(\beta)}[\mathbf{H}[y_i | \mathbf{x}_i, \mathbf{e}_i; \beta]] \\
 &\approx \mathbf{H}[\mathbf{E}_{P(\beta)}[y_i | \mathbf{x}_i, \mathbf{e}_i; \beta]] - \mathbf{E}_{P(\beta)}[\mathbf{H}[y_i | \mathbf{x}_i, \mathbf{e}_i; \beta]] \\
 &\approx -\frac{1}{T} \sum_{t=1}^T p_i^e \tilde{p}_t \cdot \log \frac{1}{T} \sum_{t=1}^T p_i^e \tilde{p}_t \\
 &\quad + \frac{1}{T} \sum_{t=1}^T p_i^e \tilde{p}_t \cdot \log p_i^e \tilde{p}_t,
 \end{aligned}$$

$$\begin{aligned}
 g(y_i, \beta | \mathbf{x}_i) &= \mathbf{H}[y_i | \mathbf{x}_i] - \mathbf{E}_{P(\beta)}[\mathbf{H}[y_i | \mathbf{x}_i; \beta]] \\
 &\approx \mathbf{H}[\mathbf{E}_{P(\beta)}[y_i | \mathbf{x}_i; \beta]] - \mathbf{E}_{P(\beta)}[\mathbf{H}[y_i | \mathbf{x}_i; \beta]] \\
 &\approx -\frac{1}{T} \sum_{t=1}^T p_t \cdot \log \frac{1}{T} \sum_{t=1}^T p_t \\
 &\quad + \frac{1}{T} \sum_{t=1}^T p_t \cdot \log p_t,
 \end{aligned}$$

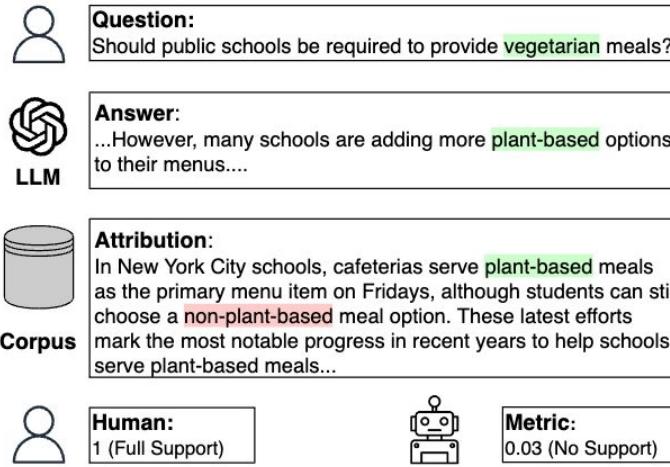
Confidence scores with and without intents

# Intent-calibrated Self-training for Answer Selection

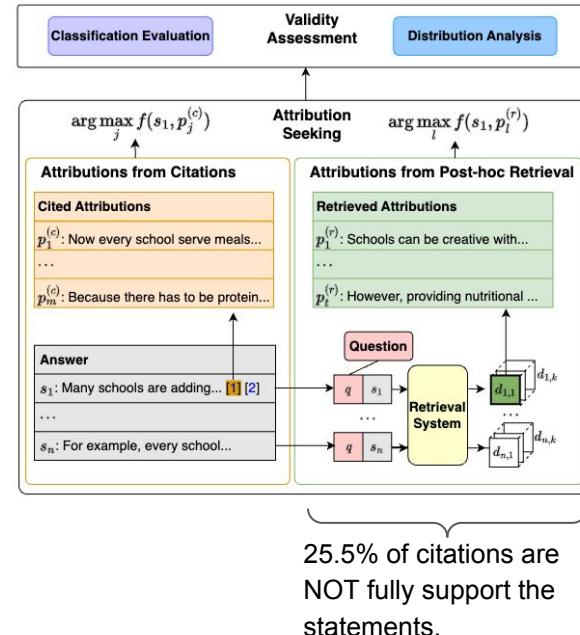
Setting	Model	MSDIALOG							MANTIS						
		P	R	F1	R@1	R@2	R@5	MAP	P	R	F1	R@1	R@2	R@5	MAP
1% labeled	IART <sup>†</sup>	22.18	46.75	30.08	25.65	46.28	77.58	47.74	48.29	52.22	50.18	50.40	68.34	86.22	66.12
	SAM <sup>†</sup>	44.17	44.36	44.26	46.89	59.06	77.02	60.72	57.75	58.62	58.18	65.10	76.32	88.54	75.60
	JM <sup>‡</sup>	44.80	44.59	44.70	44.54	60.76	84.30	61.26	62.95	62.62	62.78	62.64	77.32	92.30	75.25
	BIG <sup>‡</sup>	44.07	44.78	44.42	50.93	66.30	87.50	66.15	57.91	57.42	57.66	70.22	83.04	95.12	80.78
	GRAY <sup>‡</sup>	41.68	42.15	41.91	51.26	66.40	85.62	66.10	61.30	60.72	61.01	64.67	77.34	88.32	75.57
	GRN <sup>‡</sup>	43.41	43.37	43.39	43.28	61.60	86.46	61.19	61.75	61.10	61.42	61.06	76.64	93.66	74.56
	BERT_FP <sup>†</sup>	44.32	42.95	43.62	56.76	<b>72.08</b>	<b>91.25</b>	<b>70.90</b>	66.26	62.86	64.51	75.62	<b>86.14</b>	<b>95.22</b>	<b>84.11</b>
	BERT <sup>‡</sup>	48.56	45.34	46.90	54.79	68.32	85.80	68.04	67.28	65.62	66.44	74.82	83.00	92.16	82.41
1% labeled +all unlabeled	ICAST (Teacher)	<b>49.82</b>	<b>46.33</b>	<b>48.01</b>	<b>56.86</b>	67.81	85.38	69.03	<b>68.48</b>	<b>66.12</b>	<b>67.28</b>	<b>77.28</b>	86.12	94.98	82.98
	TSST <sup>‡</sup>	53.72	52.58	53.14	61.04	73.91	89.70	73.04	73.73	72.60	73.16	82.94	<b>91.08</b>	<b>97.88</b>	<b>89.18</b>
	ICAST	<b>57.05</b>	<b>54.32</b>	<b>55.65</b>	<b>62.21</b>	<b>76.31</b>	<b>91.07</b>	<b>73.77</b>	<b>74.89</b>	<b>72.72</b>	<b>73.79</b>	<b>83.68</b>	90.68	96.42	88.31

The predicted intents can provide more information that are useful for selecting correct answers

# Validity of Faithfulness Metrics in Attribution Seeking for LLMs

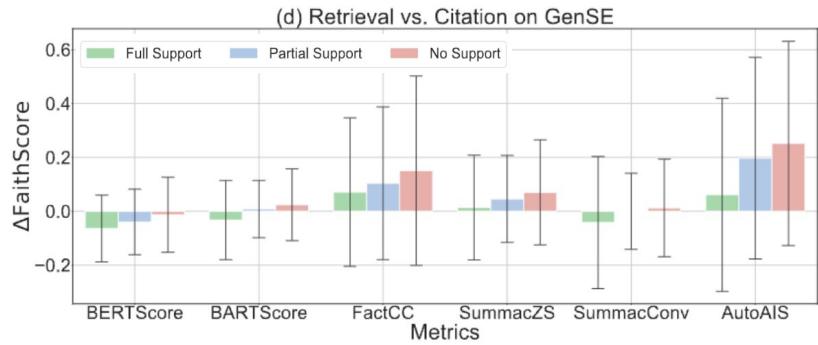


An example of an invalid faithfulness metric when the attribution contains partial support evidence.



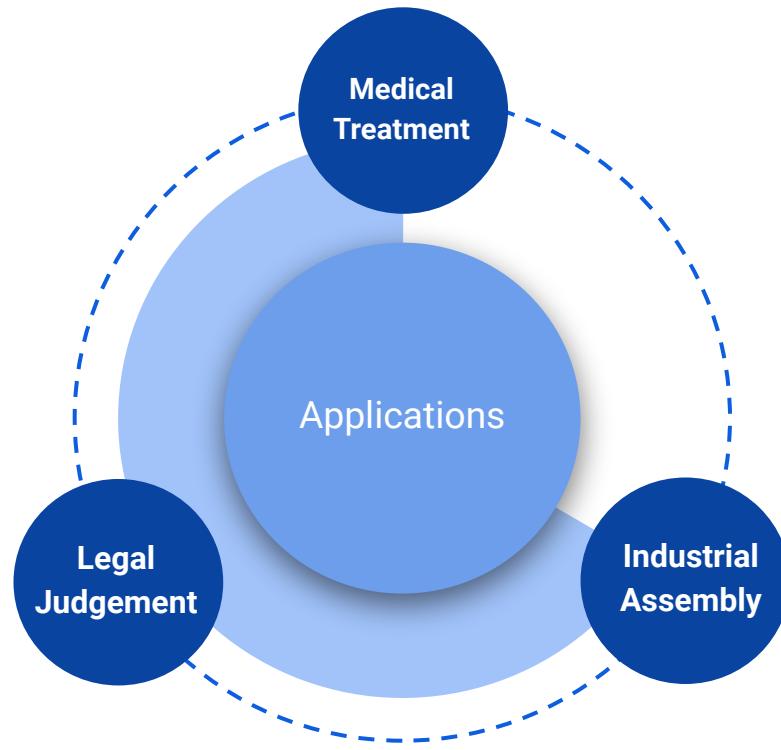
# Validity of Faithfulness Metrics in Attribution Seeking for LLMs

Metric	GENSE						ATOMIC			
	Evidence			Citation			Sentence			Fact
	F/N	F/P	P/N	F/N	F/P	P/N	F/N	F/P	P/N	F/N
BERTScore	<b>92.78</b>	79.68	<b>78.33*</b>	74.59	69.20	<b>60.50</b>	85.15	73.63	<b>68.47</b>	67.97
BARTScore	87.71	76.33	72.09	72.29	68.88	56.82	<b>85.79</b>	75.35	<b>67.82</b>	68.45
FactCC	72.09	62.56	60.74	63.16	58.79	55.08	68.83	68.41	50.99	52.56
SUMMAC <sub>ZS</sub>	89.67	79.19	<b>73.16</b>	69.22	67.83	52.92	83.24	78.96	58.83	79.32
SUMMAC <sub>Conv</sub>	71.77	78.87	35.02	70.67	<b>69.93</b>	52.32	81.70	<b>79.09</b>	56.28	<b>83.05</b>
AutoAIS	<b>92.01</b>	<b>83.46*</b>	72.20	<b>76.83</b>	<b>72.63</b>	<b>58.34</b>	<b>88.36*</b>	<b>83.59*</b>	60.33	<b>84.90</b>
Ensemble	93.93	84.64	77.08	76.84	74.11	58.05	90.33	83.88	65.72	84.03



- Partial support (F/P and P/N) drops ~10% performance compared with F/N;
- The larger the score, the better the metrics are at finding attributions that possibly contain more supporting evidence.

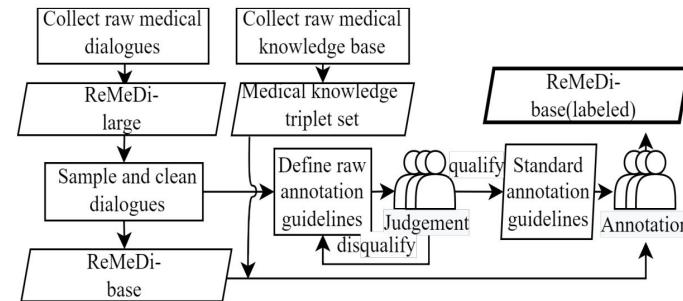
## 1.3 LLMs Tailored for Specific Applicable Domains



# Multi-Domain Multi-service Medical Dialogues

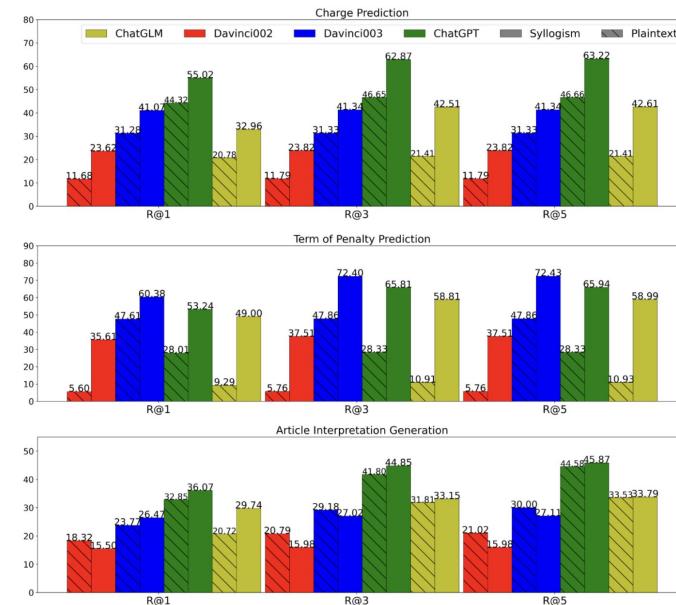


- A dataset contains 96,965 conversations between doctors and patients, including 1,557 conversations with fine-gained labels;
- Benchmarks using pretrained models (i.e., BERT-WWM, BERT-MED, GPT2, and MT5);
- A proposed self-supervised contrastive learning(SCL) model.



# Syllogistic Reasoning for Legal Judgment Analysis

<b>Input:</b> Fact Defendant A used a credit card from Bank B in May 2008 for malicious overdraft consumption of 18060 yuan without repayment.	
<b>Task 1:</b> Article Retrieval	<b>Task 2:</b> Criminal Element Generation
<b>Generated Query:</b> A324 Credit Card Fraud Crime Using Fiction of Identity or Overdraft of Credit Cards for Fraud. ... <b>Retrieved Document:</b> Article 196 [Crime of Credit Card Fraud] Whoever Maliciously overdrawn, with a relatively large amount, shall be sentenced to not more than five years or criminal detention. ... The overdraft amount reaches 18060 yuan.	<b>Subject:</b> Defendant A who is a fully capable person. <b>Subjective Element:</b> Intent <b>Object:</b> Bank B's property security <b>Objective Elements:</b> 1. Time: May 2008. 2. Location: Bank B. 3. Action: Overdraft consumption through credit card. 4. Results: The overdraft amount reaches 18060 yuan.
<b>Task 3:</b> Article Interpretation Generation	
Article 196: Defendant A maliciously overdrawn his credit card with a huge amount, which meets the Article 196 of credit card fraud.	
<b>Task 4:</b> Legal Judgment Prediction	
<b>Charge:</b> Crime of Credit Card Fraud <b>Criminal Penalty:</b> Imprisonment of up to five years	



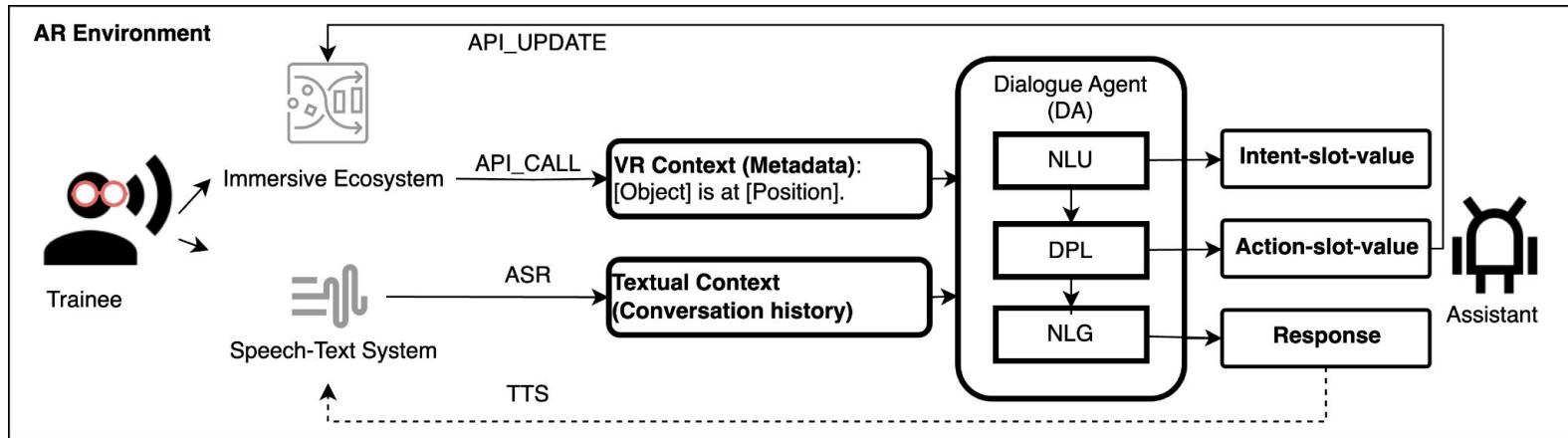
- A dataset contains 11,239 legal cases, including 1,337 facts, 124 articles, 80 charges and 5 terms.
- LLM Benchmarks: ChatGLM (6.2B+RLHF), Davinci002 (175B), Davinci003 (Davinci002+RLHF), ChatGPT(Davinci002+RLHF+Prompt&Context)

# Augmented Reality Guided Dialogue Agent for Industrial Assembly



Turn	Trainee-Assistant Conversation	Dialogue Modeling
1	Trainee: Hi, I'm ready to start assembling the car.	Trainee: StartAssembly
2	Agent: Great! Let me check the progress for you.	Agent: ProvideProgressStatus
3	Trainee: Can you please tell me <u>where I can find the engine?</u>	Trainee: GetPartLocation {PartName: Engine}
4	Agent: Sure, the engine is located in the storage area on the right side of the workshop.	Agent: ProvidePartLocation {PartName: Engine}
5	Trainee: Thank you. Now, I need <u>instructions</u> for step 5 of the assembly.	Trainee: GetAssemblyInstructions {StepNumber: 5}

# Augmented Reality Guided Dialogue Agent for Industrial Assembly



## 2. Next Steps

## 2.1 Faithfulness / Uncertainty Estimation

what is supervised tfidf learning?

Supervised TF-IDF learning is a type of supervised learning algorithm that uses the TF-IDF (Term Frequency-Inverse Document Frequency) technique for feature extraction. In this approach, the input data is represented as a bag of words, where each word is treated as a separate feature.

The TF-IDF technique is used to assign weights to each feature based on its frequency in the document and its inverse frequency across all documents. This helps to identify the most important and relevant features in the data.

In supervised TF-IDF learning, the labeled dataset is used to train a classification model that can accurately predict the class labels of new, unseen documents based on their TF-IDF feature representations. This approach is commonly used in text classification tasks, such as sentiment analysis, where the goal is to classify text documents into different categories based on their sentiment or tone.

Supervised TF-IDF learning is a powerful technique for text classification tasks, as it can effectively capture the important features of the input data and improve the accuracy of the classification model.

what is supervised tfidf learning?

About 361,000 results (0.47 seconds)

<https://www.capitalone.com> > ... > Machine Learning

**Understanding TF-IDF for Machine Learning | Capital One**  
6 Oct 2021 — TF-IDF stands for term frequency-inverse document frequency and it is a measure, used in the fields of information retrieval (IR) and machine ...

People also ask :

Is TF-IDF supervised or unsupervised?

The most popular term weighting scheme is TF-IDF (Term Frequency - Inverse Document Frequency). It is an **Unsupervised Weighting Scheme (UWS)** since it does not consider the class information in the weighting of terms. 12 Mar 2020

<https://arxiv.org> > cs

**TF-IDFC-RF: A Novel Supervised Term Weighting Scheme**

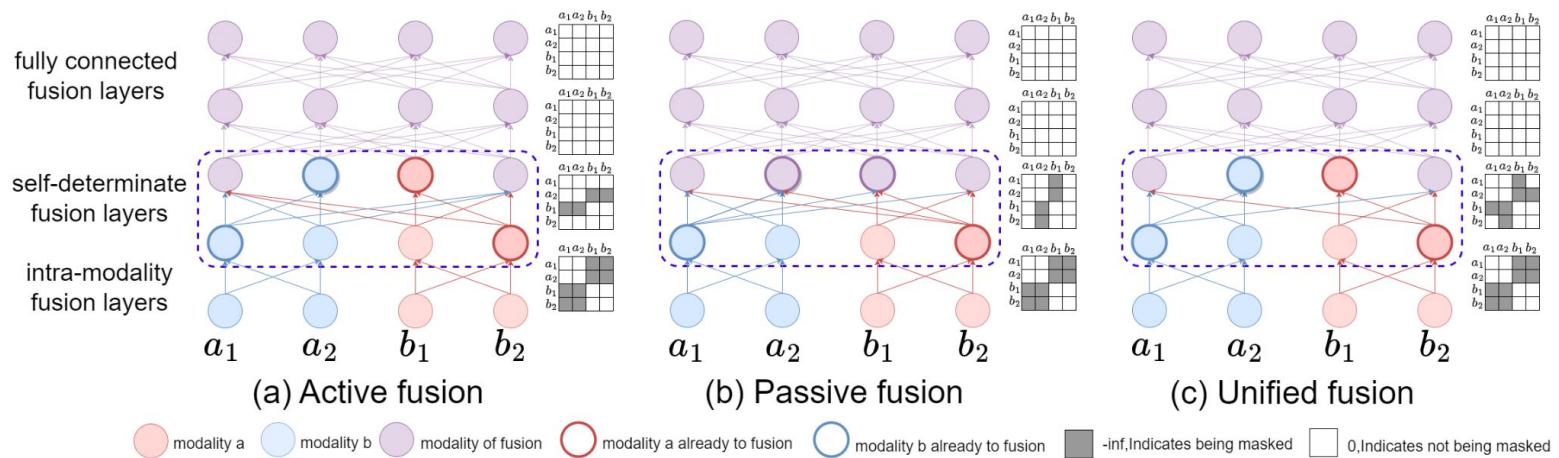
Search for: Is TF-IDF supervised or unsupervised?

What is TF-IDF in machine learning?

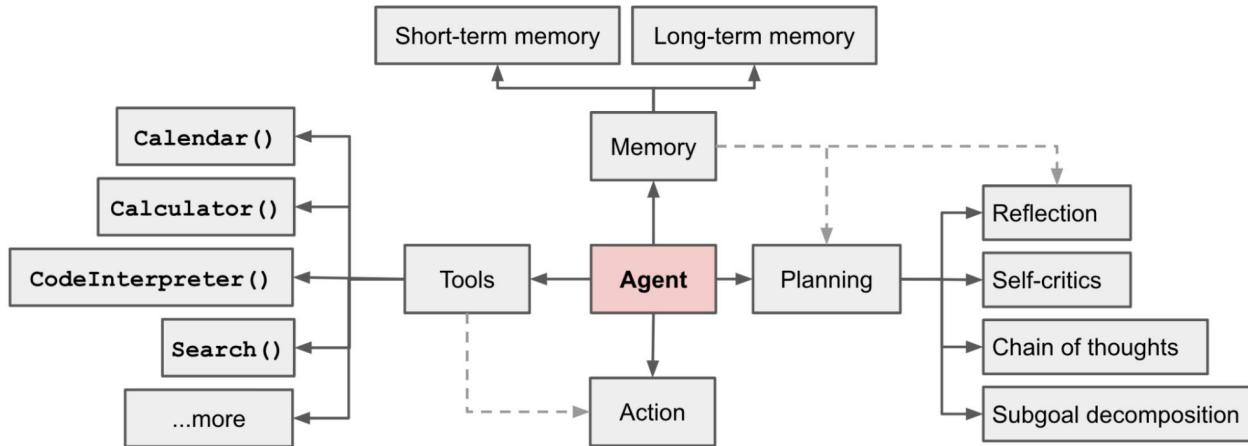
TF-IDF stands for **term frequency-inverse document frequency** and it is a measure, used in the fields of information retrieval (IR) and machine learning, that can quantify the importance or relevance of string representations (words, phrases, lemmas, etc) in a

The screenshots are from Chatgpt (left) and Google (right), respectively.

## 2.2 Multimodal LLMs



## 2.3 LLM Powered Autonomous Agents



Thank you for your attention!  
Q & A