

Knowledgeable and Multimodal Language Generation

Mohit Bansal



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

(WNGT-EMNLP 2019 Workshop)

Overall: NLG/Dialogue Model's Requirements



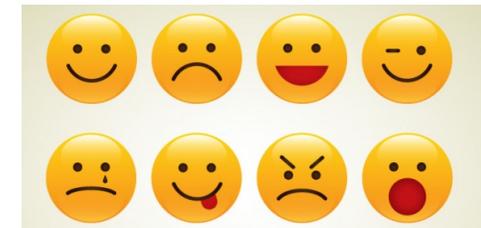
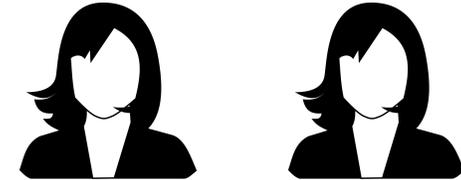
Inference in Long Context/History

Commonsense and External Knowledge

User Satisfaction Feedback & Error Robustness

Human-Personality Convincing Responses

Many-modal Grounding in Home Surroundings+Tasks (Video, Databases, etc.)



Part1: Knowledgeable and Robust NLG Models



Auxiliary Knowledge
(Entailment, Saliency)

Robustness to Missing
words, Spelling/Grammar
Errors, Paraphrases

Sensitivity to
Negations/Antonyms

Auto-Adversary
Generation

External
Commonsense

Auxiliary Knowledge via Multi-Task Learning



- MTL: Paradigm to improve generalization performance of a task using related tasks.
- The multiple tasks are learned in parallel (alternating optimization mini-batches) while using shared model representations/parameters.
- Each task benefits from extra information in the training signals of related tasks.
- Useful survey+blog by Sebastian Ruder for details of diverse MTL papers!

Auxiliary Knowledge in Language Generation



- Multi-Task & Reinforcement Learning for Entailment+Saliency Knowledge/Control in NLG (Video Captioning, Document Summarization, and Sentence Simplification)



Ground truth: A woman is slicing a red pepper.

SotA Baseline: A woman is slicing a carrot.

Our model: A woman is slicing a pepper.



Ground truth: A group of boys are fighting.

SotA Baseline: A group of men are dancing.

Our model: Two men are fighting.

Document: *top activists arrested after last month 's anti-government rioting are in good condition , a red cross official said saturday .*

Ground-truth: *arrested activists in good condition says red cross*

SotA Baseline: *red cross says it is good condition after riots*

Our model: *red cross says detained activists in good condition*

Document: *canada 's prime minister has dined on seal meat in a gesture of support for the sealing industry .*

Ground-truth: *canadian pm has seal meat*

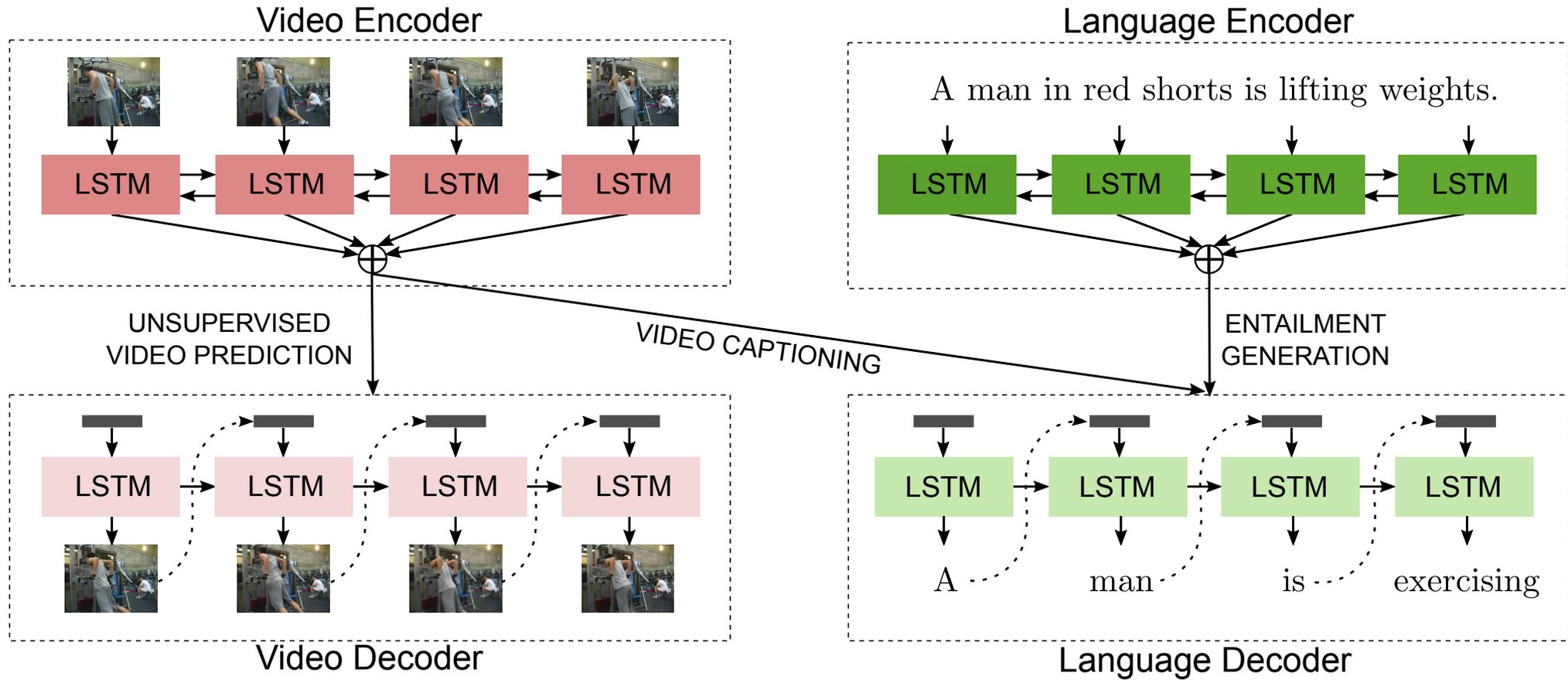
SotA Baseline: *canadian pm says seal meat is a matter of support*

Our model: *canada 's prime minister dines with seal meat*

Auxiliary Knowledge in Language Generation



- Many-to-Many Multi-Task Learning for Video Captioning (with Video and Entailment Generation)



Results (YouTube2Text)



Models	METEOR	CIDEr-D	ROUGE-L	BLEU-4
PREVIOUS WORK				
LSTM-YT (Venugopalan et al., 2015b)	26.9	-	-	31.2
S2VT (Venugopalan et al., 2015a)	29.8	-	-	-
Temporal Attention (Yao et al., 2015)	29.6	51.7	-	41.9
LSTM-E (Pan et al., 2016b)	31.0	-	-	45.3
Glove + DeepFusion (Venugopalan et al., 2016)	31.4	-	-	42.1
p-RNN (Yu et al., 2016)	32.6	65.8	-	49.9
HNRE + Attention (Pan et al., 2016a)	33.9	-	-	46.7
OUR BASELINES				
Baseline (V)	31.4	63.9	68.0	43.6
Baseline (G)	31.7	64.8	68.6	44.1
Baseline (I)	33.3	75.6	69.7	46.3
Baseline + Attention (V)	32.6	72.2	69.0	47.5
Baseline + Attention (G)	33.0	69.4	68.3	44.9
Baseline + Attention (I)	33.8	77.2	70.3	49.9
Baseline + Attention (I) (E) \otimes	35.0	84.4	71.5	52.6
OUR MULTI-TASK LEARNING MODELS				
\otimes + Video Prediction (1-to-M)	35.6	88.1	72.9	54.1
\otimes + Entailment Generation (M-to-1)	35.9	88.0	72.7	54.4
\otimes + Video Prediction + Entailment Gener (M-to-M)	36.0	92.4	72.8	54.5

* All models (1-to-M, M-to-1 and M-to-M) stat. signif. better than strong SotA baseline.

Results (MSR-VTT)



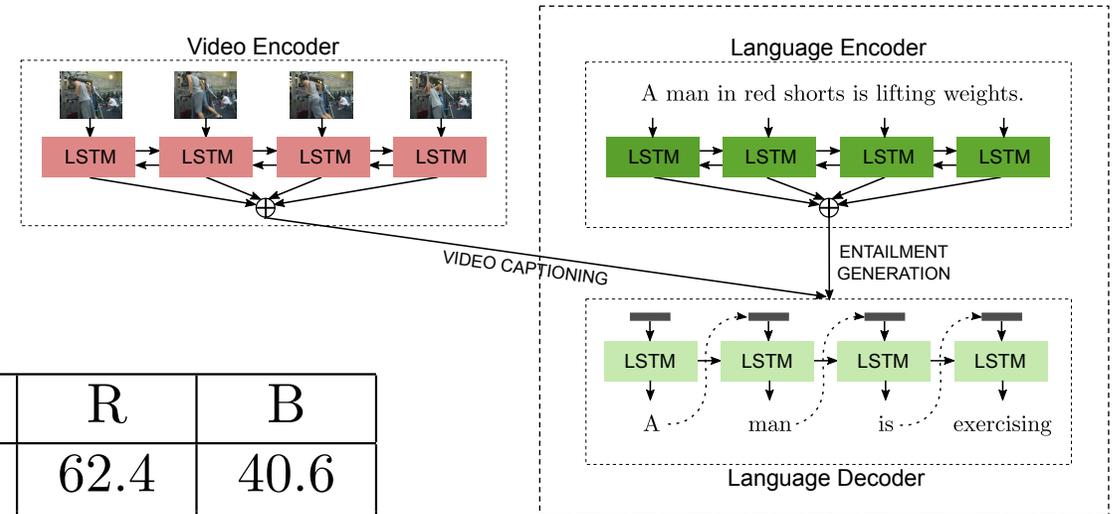
- Diverse video clips from a commercial video search engine

Models	METEOR	CIDEr-D	ROUGE-L	BLEU-4
Venugopalan et al., 2015	23.4	-	-	32.3
Yao et al., 2015	25.2	-	-	35.2
Xu et al., 2016	25.9	-	-	36.6
Rank1: v2t_navigator	28.2	44.8	60.9	40.8
Rank2: Aalto	26.9	45.7	59.8	39.8
Rank3: VideoLAB	27.7	44.1	60.6	39.1
Our Model (New Rank1)	28.8	47.1	60.2	40.8

Results (Entailment Generation)



- Video captioning mutually also helps improve the entailment-generation task in turn (w/ statistical significance)



Models	M	C	R	B
Entailment Generation	29.6	117.8	62.4	40.6
+Video Caption (M-to-1)	30.0	121.6	63.9	41.6

- New multi-reference split setup of SNLI to allow automatic metric evaluation and a zero train-test premise overlap

Human Evaluation



- Multi-task model > strong non-multitask baseline on relevance and coherence/fluency (for both video captioning and entailment generation)

	YouTube2Text		Entailment	
	Relev.	Coher.	Relev.	Coher.
Not Distinguish.	70.7%	92.6%	84.6%	98.3%
SotA Baseline Wins	12.3%	1.7%	6.7%	0.7%
Multi-Task Wins	17.0%	5.7%	8.7%	1.0%

Analysis Examples



Ground truth: Two women are shopping in a store.
Two girls are shopping.

Baseline model: A man is doing a monkey in a store.

Multi-task model: A woman is shopping in a store.



Ground truth: Two men are fighting.
A group of boys are fighting.

Baseline model: A group of men are dancing.

Multi-task model: Two men are fighting.

(a) complex examples where the multi-task model performs better than baseline

Analysis Examples



Ground truth: A woman slices a shrimp tail.
A girl is cutting a fish tale.
Baseline model: A person is cutting the something.
Multi-task model: A woman is cutting a piece of meat.



Ground truth: Two men are talking aggressively.
The boy is talking.
Baseline model: A man is crying.
Multi-task model: A man is talking.

(b) ambiguous examples (i.e., ground truth itself confusing) where multi-task model still correctly predicts one of the possible categories

Analysis Examples



Ground truth: A monkey and a deer are fighting.
A gazelle is fighting with a baboon.
Baseline model: A man is walking on the ground.
Multi-task model: A monkey is walking.



Ground truth: A dog climbs into a dryer.
A dog is in a washing machine.
Baseline model: A man is playing.
Multi-task model: A man is playing with a toy.

(c) complex examples where both models perform poorly

(d) baseline > MTL: both correct but low specificity

- Overall, multi-task model's captions are better at both temporal action prediction and logical entailment w.r.t. ground truth captions (ablated examples in paper).

Auxiliary Knowledge in Language Generation



- Reverse Multi-Task Benefits: Improved Entailment Generation

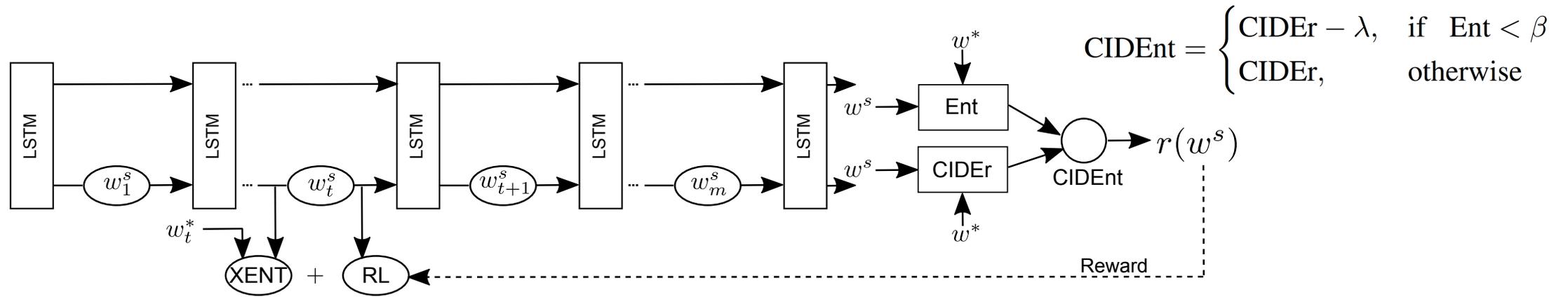
Given Premise	Generated Entailment
a man on stilts is playing a tuba for money on the boardwalk	a man is playing an instrument
a child that is dressed as spiderman is ringing the doorbell	a child is dressed as a superhero
several young people sit at a table playing poker	people are playing a game
a woman in a dress with two children	a woman is wearing a dress
a blue and silver monster truck making a huge jump over crushed cars	a truck is being driven

Auxiliary Knowledge in Language Generation



- RL Reward = Entailment-corrected phrase-matching metrics such as CIDEr \rightarrow CIDEnt

Ground-truth caption	Generated (sampled) caption	CIDEr	Ent
a man is spreading some butter in a pan	puppies is melting butter on the pan	140.5	0.07
a panda is eating some bamboo	a panda is eating some fried	256.8	0.14
a monkey pulls a dogs tail	a monkey pulls a woman	116.4	0.04
a man is cutting the meat	a man is cutting meat into potato	114.3	0.08
the dog is jumping in the snow	a dog is jumping in cucumbers	126.2	0.03
a man and a woman is swimming in the pool	a man and a whale are swimming in a pool	192.5	0.02



Auxiliary Knowledge in Language Generation



Models	BLEU-4	METEOR	ROUGE-L	CIDEr-D	CIDEnt	Human*
PREVIOUS WORK						
Venugopalan (2015b)*	32.3	23.4	-	-	-	-
Yao et al. (2015)*	35.2	25.2	-	-	-	-
Xu et al. (2016)	36.6	25.9	-	-	-	-
Pasunuru and Bansal (2017)	40.8	28.8	60.2	47.1	-	-
Rank1: v2t_navigator	40.8	28.2	60.9	44.8	-	-
Rank2: Aalto	39.8	26.9	59.8	45.7	-	-
Rank3: VideoLAB	39.1	27.7	60.6	44.1	-	-
OUR MODELS						
Cross-Entropy (Baseline-XE)	38.6	27.7	59.5	44.6	34.4	-
CIDEr-RL	39.1	28.2	60.9	51.0	37.4	11.6
CIDEnt-RL (New Rank1)	40.5	28.4	61.4	51.7	44.0	18.4

Table 2: Our primary video captioning results on MSR-VTT (CIDEnt-RL is stat. significantly better than CIDEr-RL in all metrics, and CIDEr-RL is better than Baseline-XE).

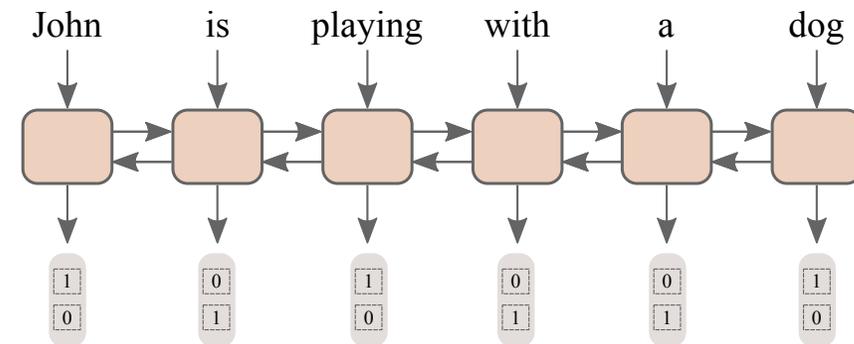
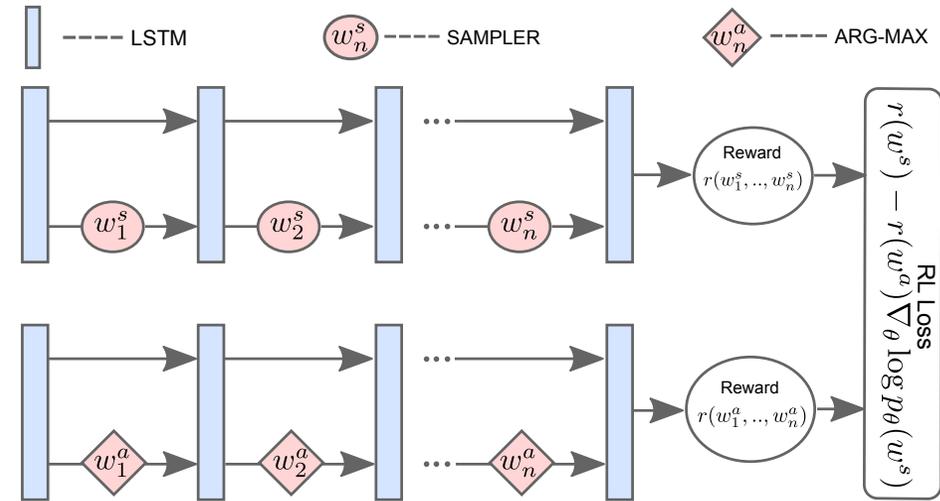
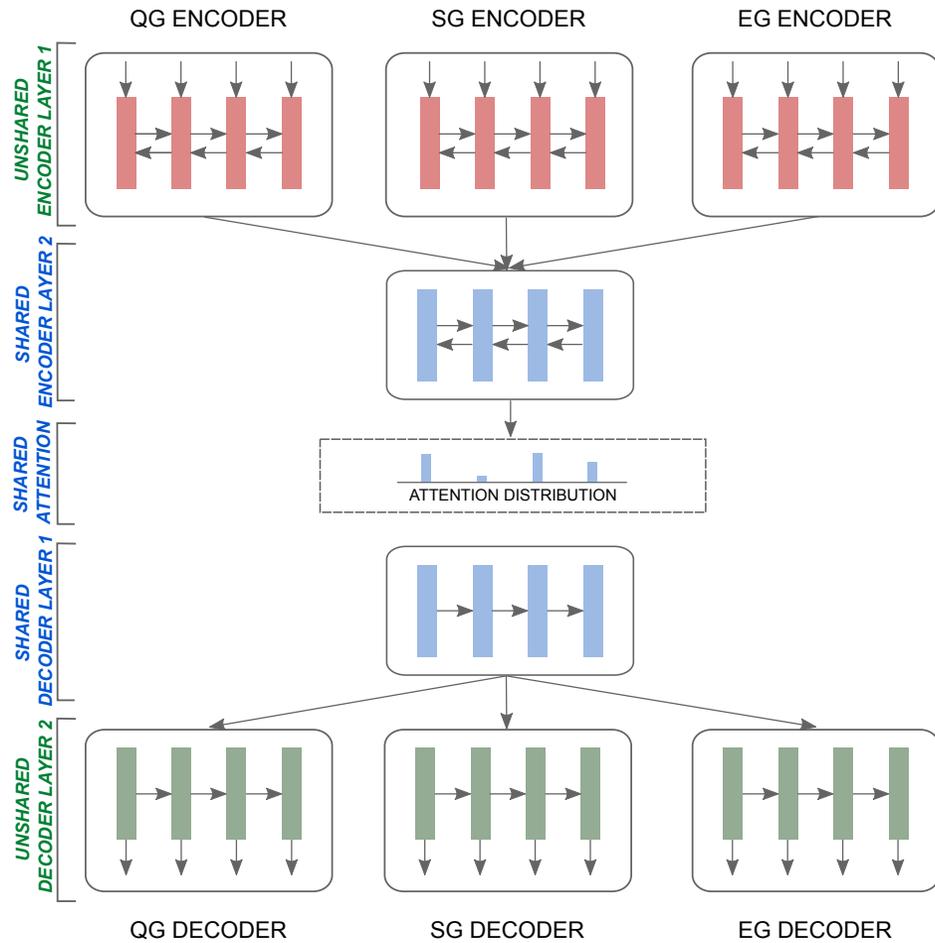
	Relevance	Coherence		Relevance	Coherence
Not Distinguishable	64.8%	92.8%	Not Distinguishable	70.0%	94.6%
Baseline-XE Wins	13.6%	4.0%	CIDEr-RL Wins	11.6%	2.8%
CIDEr-RL Wins	21.6%	3.2%	CIDEnt-RL Wins	18.4%	2.8%

Table 3: Human evaluation results on MSR-VTT (CIDEnt-RL is stat. significantly better than CIDEr-RL, and CIDEr-RL is better than Baseline-XE).

Auxiliary Knowledge in Language Generation



- Multi-Task & Reinforcement Learning with Entailment+Saliency Knowledge for Summarization



Auxiliary Knowledge in Language Generation



Input Document: celtic have written to the scottish football association in order to gain an ‘understanding’ of the refereeing decisions during their scottish cup semi-final defeat by inverness on sunday . the hoops were left outraged by referee steven mclean ’s failure to award a penalty or red card for a clear handball in the box by josh meekings to deny leigh griffith ’s goal-bound shot during the first-half . caley thistle went on to win the game 3-2 after extra-time and denied rory delia ’s men the chance to secure a domestic treble this season . celtic striker leigh griffiths has a goal-bound shot blocked by the outstretched arm of josh meekings after the restart for scything down marley watkins in the area . greg tansey duly converted the resulting penalty . edward ofere then put caley thistle ahead , only for john guidetti to draw level for the bhoys . with the game seemingly heading for penalties , david raven scored the winner on 117 minutes , breaking thousands of celtic hearts . celtic captain scott brown -lrb- left -rrb- protests to referee steven mclean but the handball goes unpunished . griffiths shows off his acrobatic skills during celtic ’s eventual surprise defeat by inverness . celtic pair aleksandar tonev -lrb- left -rrb- and john guidetti look dejected as their hopes of a domestic treble end .

Ground-truth Summary: celtic were defeated 3-2 after extra-time in the scottish cup semi-final . leigh griffiths had a goal-bound shot blocked by a clear handball. however, no action was taken against offender josh meekings. the hoops have written the sfa for an ‘understanding’ of the decision .

See et al. (2017): **john hartson** was once on the end of a major **hampden injustice** while playing for celtic . but he can not see any point in his old club writing to the scottish football association over the latest controversy at the national stadium . hartson had a goal wrongly disallowed for offside while celtic were leading 1-0 at the time but went on to lose 3-2 .

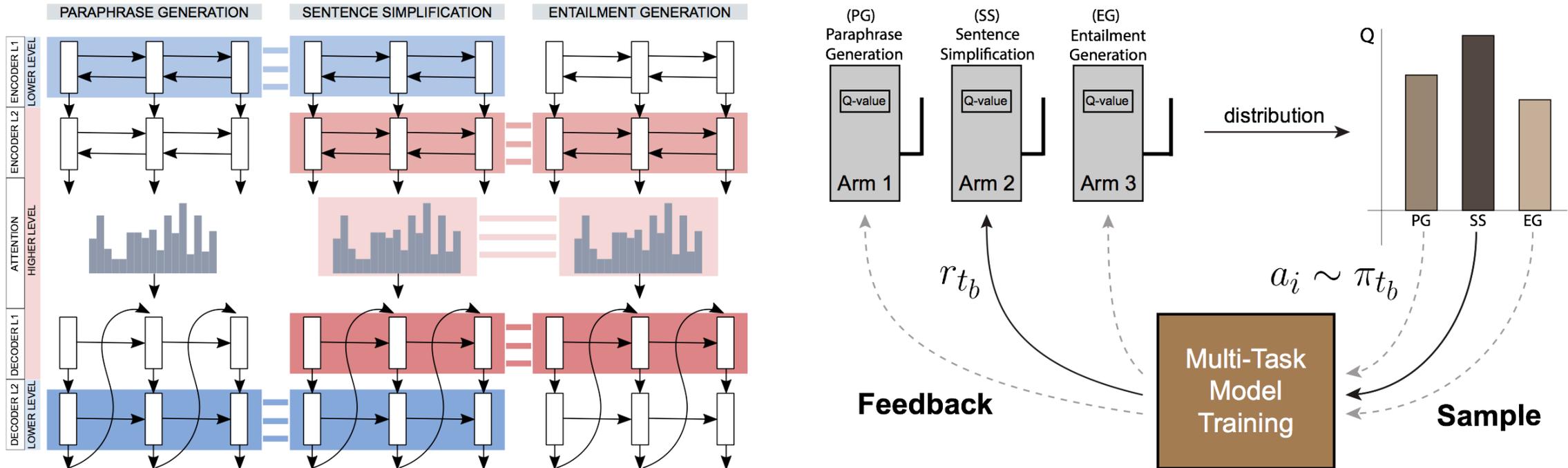
Our Baseline: **john hartson** scored the late winner in 3-2 win against celtic . celtic were leading 1-0 at the time but went on to lose 3-2 . some fans have questioned how referee steven mclean and **additional assistant alan muir** could have missed the infringement .

Our Multi-task Summary: celtic have written to the scottish football association in order to gain an ‘ understanding ’ of the refereeing decisions . the hoops were left outraged by referee steven mclean ’s failure to award a penalty or red card for a clear handball in the box by josh meekings . celtic striker leigh griffiths has a goal-bound shot blocked by the outstretched arm of josh meekings .

Auxiliary Knowledge in Language Generation

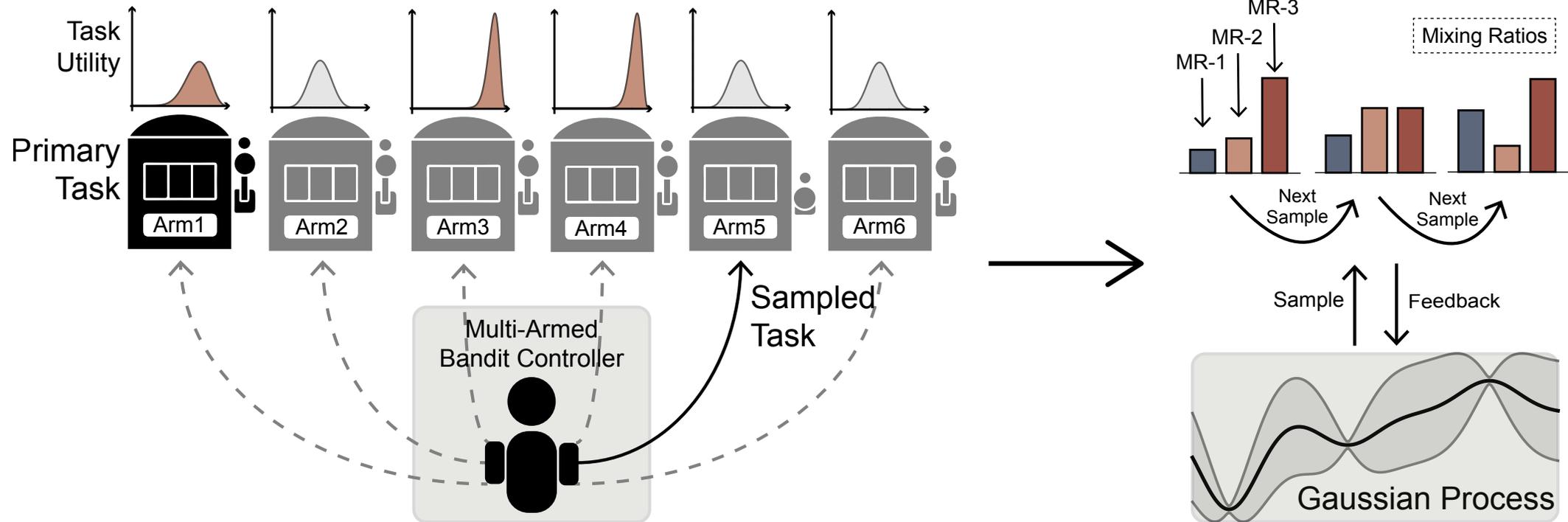


- Dynamic-Curriculum MTL with Entailment+Paraphrase Knowledge for Sentence Simplification



Code: <https://github.com/HanGuo97/MultitaskSimplification>

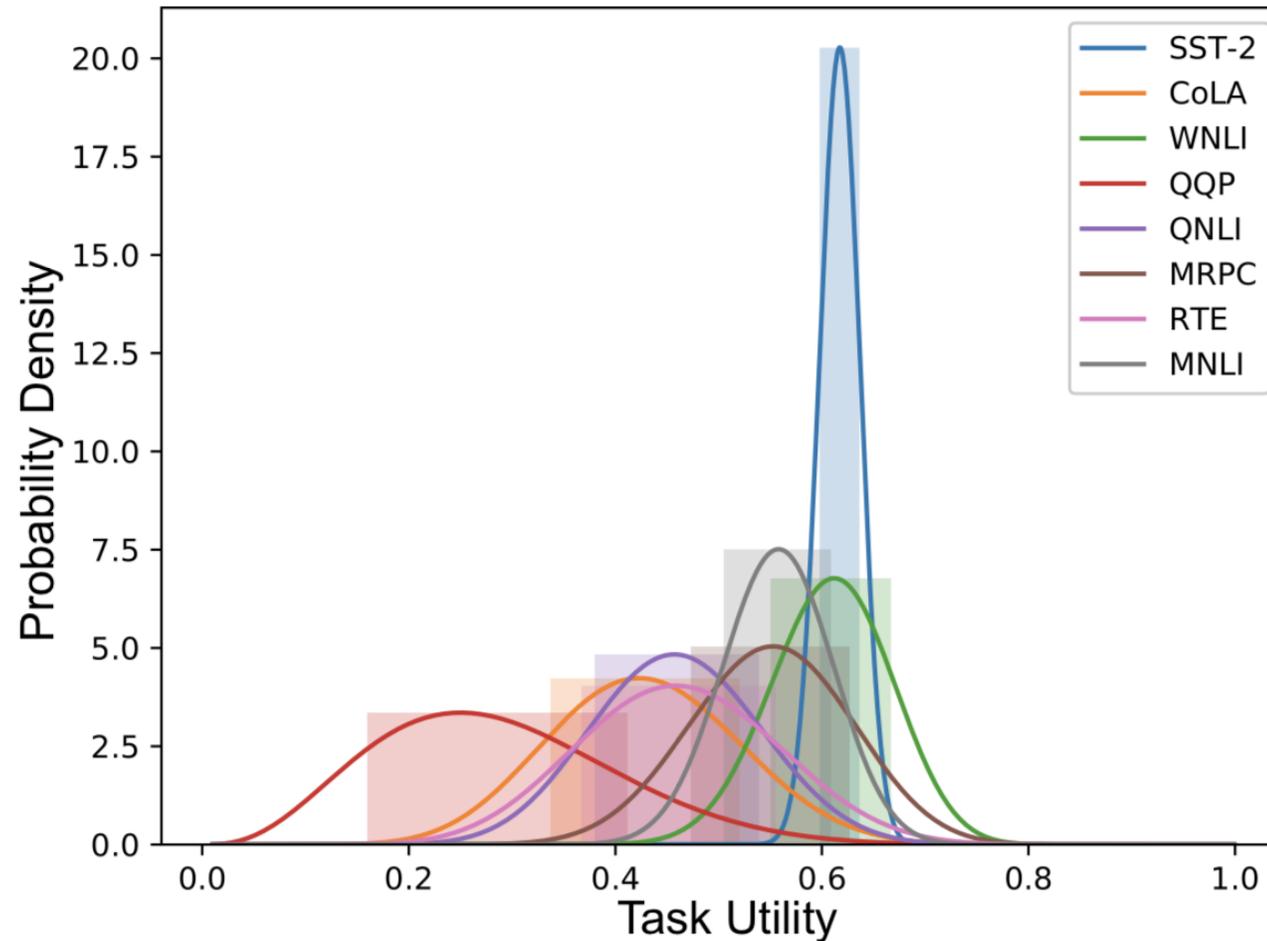
AutoSeM: Automatic Auxiliary Task Selection+Mixing



Left: the multi-armed bandit controller used for task selection, where each arm represents a candidate auxiliary task. The agent iteratively pulls an arm, observes a reward, updates its estimates of the arm parameters, and samples the next arm. Right: the Gaussian Process controller used for automatic mixing ratio (MR) learning. The GP controller sequentially makes a choice of mixing ratio, observes a reward, updates its estimates, and selects the next mixing ratio to try, based on the full history of past observations.

Code: <https://github.com/HanGuo97/AutoSeM>

Interpretability: Visualization of Stage-1 Task Selection

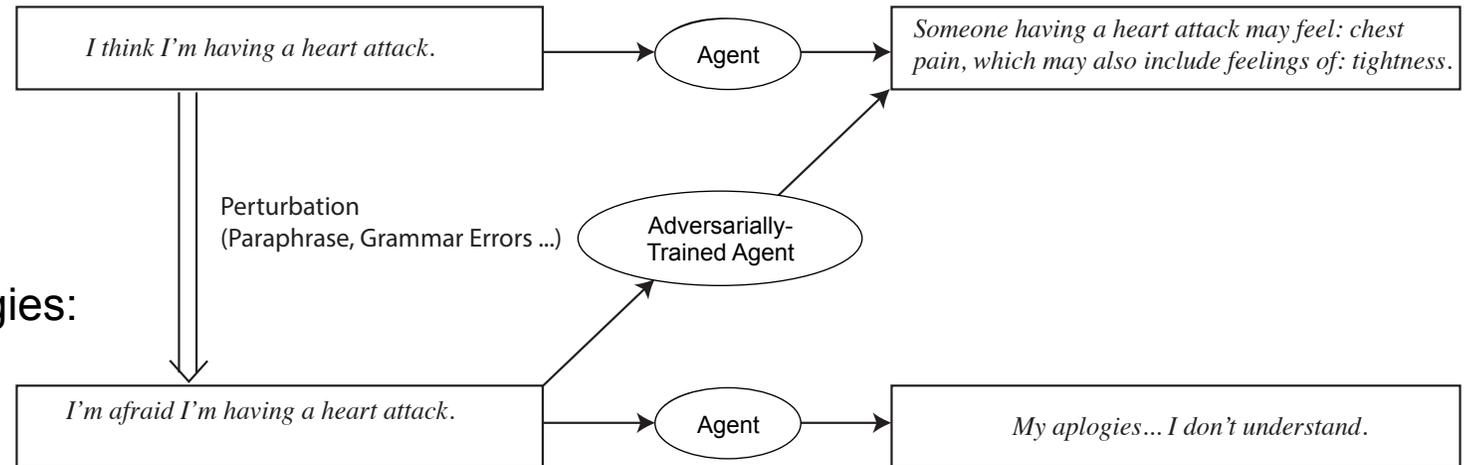


Visualization of task utility estimates from the multi-armed bandit controller on SST-2 (primary task). The x-axis represents the task utility, and the y-axis represents the corresponding probability density. Each curve corresponds to a task and the bar corresponds to their confidence interval.

Adversarially-Robust Dialogue Generation



- Robustness to real-world noise (e.g., user errors) and subtle but important markers!
- “Should-Not-Change” Over-Sensitivity Strategies:
 - Random Swap
 - Stopword Dropout
 - Data-level Paraphrasing
 - Generative-level Paraphrasing
 - Grammar Errors
- “Should-Change” Over-Stability Strategies:
 - Add Negation
 - Antonym
 - Random Inputs
 - Random Inputs with Preserved Entities
 - Confusing Entity
- Tasks/Datasets: Ubuntu (Activity/Entity F1, Human Eval), CoCoA (Completion Rate)
- Models: VHRED, Reranking-RL, DynoNet



Adversarially-Robust Dialogue Generation



Strategy Name	N-train + A-test	A-train + A-test	A-train + N-test	N-train + N-test
Normal Input	-	-	-	5.94, 3.52
Random Swap	6.10, 3.42	6.47, 3.64	6.42, 3.74	-
Stopword Dropout	5.49, 3.44	6.23, 3.82	6.29, 3.71	-
Data-Level Para.	5.38, 3.18	6.39, 3.83	6.32, 3.87	-
Generative-Level Para.	4.25, 2.48	5.89, 3.60	6.11, 3.66	-
Grammar Errors	5.60, 3.09	5.93, 3.67	6.05, 3.69	-
All Should-Not-Change	-	-	6.74, 3.97	-
Add Negation	6.06, 3.42	5.01, 3.12	6.07, 3.46	-
Antonym	5.85, 3.56	5.43, 3.43	5.98, 3.56	-

Table 2: Activity and Entity F1 results of adversarial strategies on the **VHRED** model.

	VHRED	Tie	Combined-VHRED
Winning %	28	22	49

Table 5: Human evaluation results on comparison between VHRED and VHRED train on all Should-Not-Change strategies combined.

Adversarially-Robust Dialogue Generation



Context	Response
<p>N: ... you could save your ubuntu files and reinstall Windows , then install ubuntu as a dual boot option __eou__ __eot__ aight buddy , so how do i get that **unknown** space back__eou__</p> <p>Random Swap: ... you could your save ubuntu and files Windows reinstall , then install ubuntu as dual a option boot __eou__ __eot__ aight buddy , so do how i that get space **unknown** back __eou__</p>	<p>NN: you can use the Live CD , you can install Ubuntu on the same partition as the Windows partition __eou__</p> <p>NA: I am using ubuntu . __eou__</p> <p>AA: you can use Windows XP on the Windows partition , and then install Ubuntu on the same drive __eou__</p>
<p>N: ... ok then the problem is you don't have the nvidia driver __eou__ dpkg -l — grep nvidia __eou__ __eot__ i can see some nvidia thing :P __eou__</p> <p>Stopword Dropout: ... ok then problem is you don't have nvidia driver __eou__ dpkg -l — grep nvidia __eou__ __eot__ i see nvidia thing :P</p>	<p>NN: sudo apt-get install nvidia-glx __eou__</p> <p>NA: ok __eou__</p> <p>AA: you can use the nvidia driver __eou__</p>

Auto-Augment Adversary Generation



Ribeiro et al., 2018; Zhao et al., 2018

How do we automatically generate the best adversaries without manual design? Our AutoAugment model consists of a controller and a target model. The controller first samples a policy that transforms the original data to augmented data, on which the target model trains. After training, the target model is evaluated to obtain the performance on the validation set. This performance is then fed back to the controller as the reward signal.

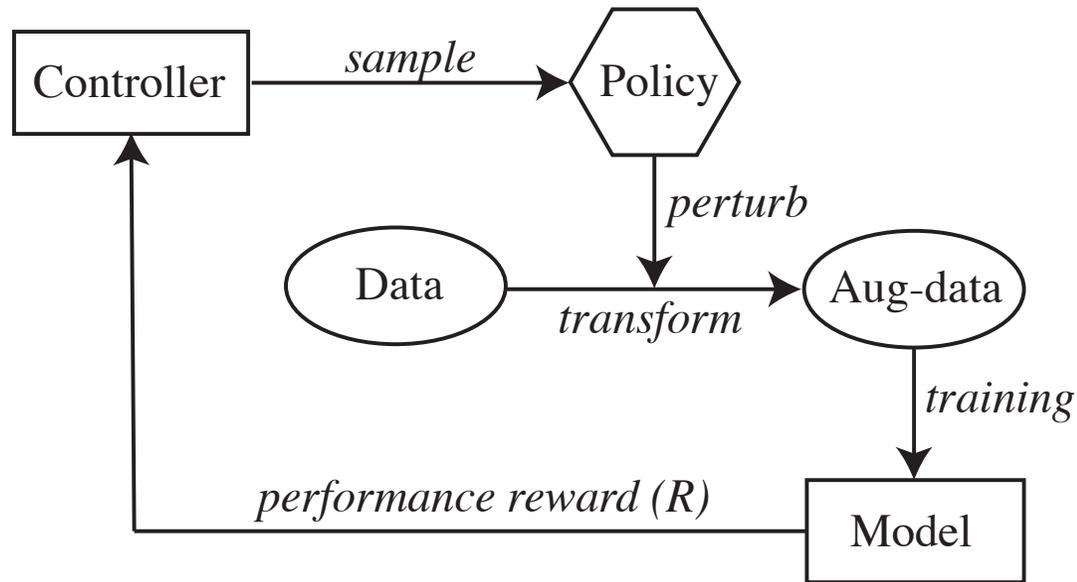


Figure 1: The controller samples a policy to perturb the training data. After training on the augmented inputs, the model feeds the performance back as reward.

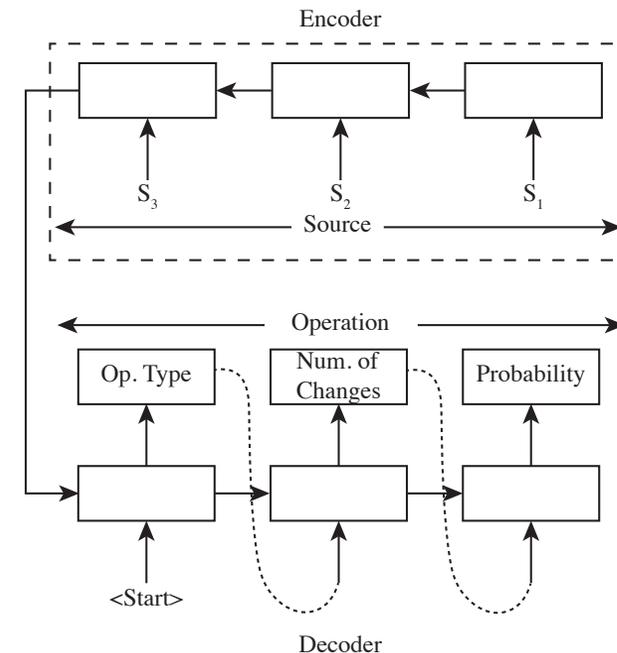


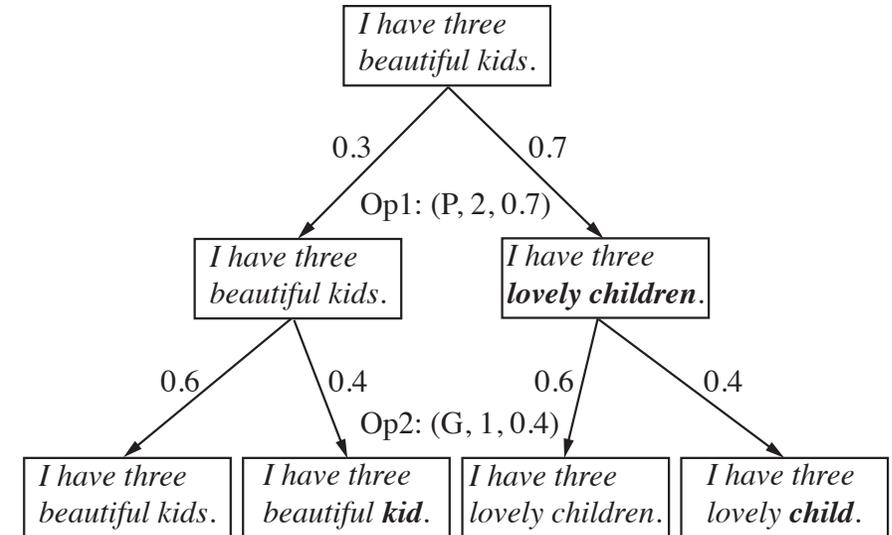
Figure 3: AutoAugment controller. An input-agnostic controller corresponds to the lower part of the figure. It samples a list of operations in sequence. An input-aware controller additionally has an encoder (upper part) that takes in the source inputs of the data.

Auto-Augment Adversary Generation



Policy Hierarchy and Search Space:

- A policy consists of 4 sub-policies;
- Each sub-policy consists of 2 operations applied in sequence;
- Each operation is defined by 3 parameters: **Operation Type**, **Number of Changes** (the maximum number of times allowed to perform the operation, and the **Probability** of applying that operation.
- Our pool of operations contains **Random Swap**, **Stopword Dropout**, **Paraphrase**, **Grammar Errors**, and **Stammer**.



Subdivision of Operations:

- **Stopword Dropout:** To allow the controller to learn more nuanced combinations of operations, divide Stopword Dropout into 7 categories: Noun, Adposition, Pronoun, Adverb, Verb, Determiner, and Other.
- **Grammar Errors:** Noun (plural/singular confusion) and Verb (verb inflected/base form confusion).

Figure 2: Example of a sub-policy applied to a source input. E.g., the first operation (Paraphrase, 2, 0.7) paraphrases the input twice with probability 0.7.

Auto-Augment Adversary Generation



- **Setup:** Variational Hierarchical Encoder-Decoder (VHRED) (Serban et al., 2017b) on troubleshooting Ubuntu Dialogue task (Lowe et al., 2015); REINFORCE (Williams, 1992; Sutton et al., 2000) to train the controller.
- **Evaluation:** Serban et al. (2017a), evaluate on F1s for both activities (technical verbs) and entities (technical nouns). We also conducted human studies on Mturk, comparing each of the input-agnostic/aware models with the VHRED baseline and All-operations from Niu and Bansal (2018).

	Activity F1	Entity F1
LSTM	1.18	0.87
HRED	4.34	2.22
VHRED	4.63	2.53
VHRED (w/ attn.)	5.94	3.52
All-operations	6.53	3.79
Input-aware	7.04	3.90
Input-agnostic	7.02	4.00

Table 1: Activity, Entity F1 results reported by previous work, the All-operations and AutoAugment models.

	W	T	L	W - L
Input-agnostic vs. baseline	48	23	29	19
Input-aware vs. baseline	45	27	28	17
Input-agnostic vs. All-ops	43	27	30	13
Input-aware vs. All-ops	50	13	37	13

Table 4: Top 3 policies on the validation set and their test performances. Operations: R=Random Swap, D=Stopword Dropout, P=Paraphrase, G=Grammar Errors, S=Stammer. Universal tags: n=noun, v=verb, p=pronoun, adv=adverb, adp=adposition.

Sub-policy1	Sub-policy2	Sub-policy3	Sub-policy4
P, 1, 0.5	D _v , 3, 0.2	R, 3, 0.9	D _p , 2, 0.3
D _{adv} , 4, 0.4	R, 1, 0.5	D _{adp} , 1, 0.5	D _{adp} , 2, 0.1
D _n , 1, 0.8	D _o , 3, 1.0	P, 4, 0.4	G _n , 3, 0.3
G _v , 1, 0.9	D _o , 3, 0.1	S, 3, 0.4	R, 1, 0.2
D _v , 2, 0.5	D _v , 2, 0.7	S, 3, 0.5	P, 1, 1.0
R, 2, 0.2	G _v , 1, 0.9	D _o , 1, 0.5	G _n , 2, 0.6

Table 2: Human evaluation results on comparisons among the baseline, All-operations, and the two AutoAugment models. W: Win, T: Tie, L: Loss.



Auto-Augment Adversary Generation

- **Setup:** Variational Hierarchical Encoder-Decoder (VHRED) (Serban et al., 2017b) on troubleshooting Ubuntu Dialogue task (Lowe et al., 2015); REINFORCE (Williams, 1992; Sutton et al., 2000) to train the controller.
- **Evaluation:** Serban et al. (2017a), evaluate on F1s for both activities (technical verbs) and entities (technical nouns). We also conducted human studies on Mturk, comparing each of the input-agnostic/aware models with the VHRED baseline and All-operations from Niu and Bansal (2018).

	Activity F1	Entity F1
LSTM	1.18	0.87
HRED	4.34	2.22
VHRED	4.63	2.52
VHRED (w/ attn.)	5.94	3.51
All-operations	6.53	3.71
Input-aware	7.04	3.91
Input-agnostic	7.02	4.01

Table 1: Activity, Entity F1 results reported by previous work, the All-operations and AutoAugment model.

	W	T	L	W - L
Input-agnostic vs. baseline	48	23	29	19
Input-aware vs. baseline	45	27	28	17
Input-agnostic vs. All-ops	43	27	30	13
			7	13

performances. Operations: Grammar Errors, =adverb, adp=adposition.

Still several challenges: better AutoAugm algorithms for RL speed, reward sparsity, other NLU/NLG tasks? Visit Tong's poster Nov5 3.30pm for more details!

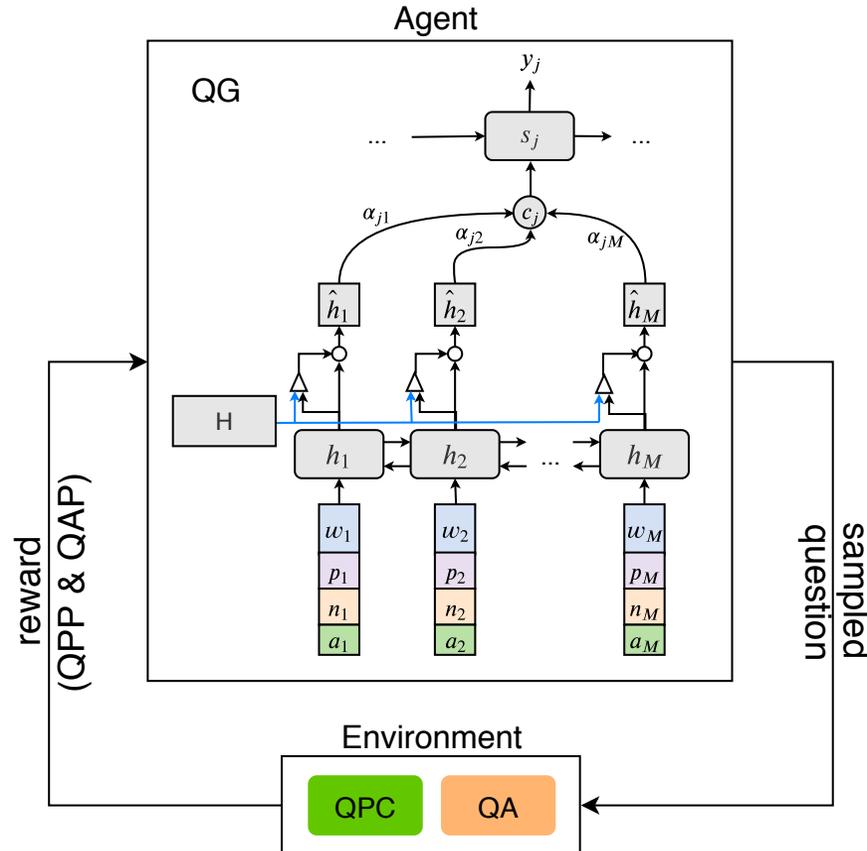
$D_{adv}, 4, 0.4$	$R, 1, 0.5$	$D_{adp}, 1, 0.5$	$D_{adp}, 2, 0.1$
$D_n, 1, 0.8$	$D_o, 3, 1.0$	$P, 4, 0.4$	$G_n, 3, 0.3$
$G_v, 1, 0.9$	$D_o, 3, 0.1$	$S, 3, 0.4$	$R, 1, 0.2$
$D_v, 2, 0.5$	$D_v, 2, 0.7$	$S, 3, 0.5$	$P, 1, 1.0$
$R, 2, 0.2$	$G_v, 1, 0.9$	$D_o, 1, 0.5$	$G_n, 2, 0.6$

Table 2: Human evaluation results on comparisons among the baseline, All-operations, and the two AutoAugment models. W: Win, T: Tie, L: Loss.

Question Generation with Semantic Validity Knowledge



- “Semantic drift” problem
 - Generated questions semantically drift away from the given context and answer .



Context: ...during the age of enlightenment, philosophers such as **john locke** advocated the principle in their writings, whereas others, such as thomas hobbes, strongly opposed it. montesquieu was one of the foremost supporters of separating the legislature, the executive, and the judiciary...

Gt: who was an advocate of separation of powers?

Base: who opposed the principle of enlightenment?

Ours: who advocated the principle in the age of enlightenment?

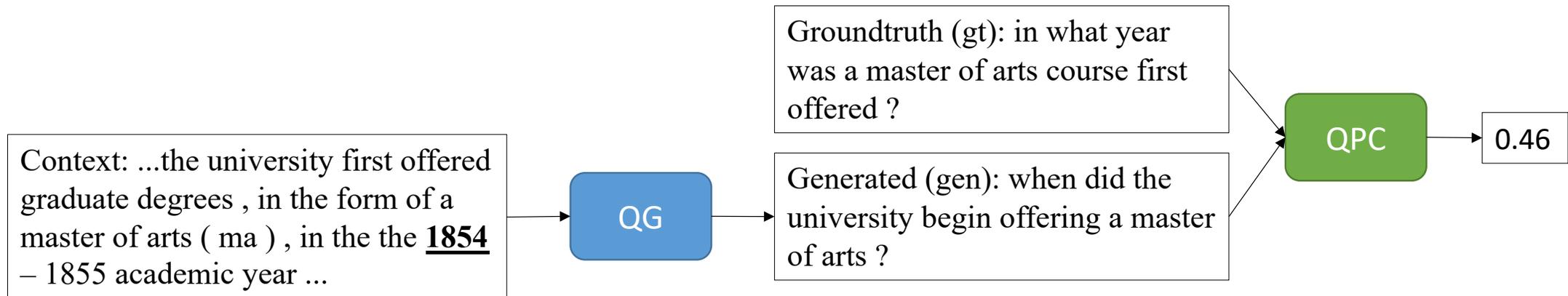
- Two “semantics-enhanced” rewards
 - QPP: Question Paraphrasing Probability
 - QAP: Question Answering Probability
- Reinforcement learning:
 - Policy gradient (Williams, 1992)
 - Mixed loss (Paulus et al., 2017)
 - Multi-reward optimization (Pasunuru & Bansal, 2018)

Question Generation with Semantic Validity Knowledge



- QPP (**Q**uestion **P**araphrasing **P**robability) reward:
 - From QPC (**Q**uestion **P**araphrasing **C**lassification) model
 - Represents “the probability of the generated question and the ground-truth question being paraphrases”

$$p_{qpc}(is_para = true | q_{gt}, q_{gen})$$

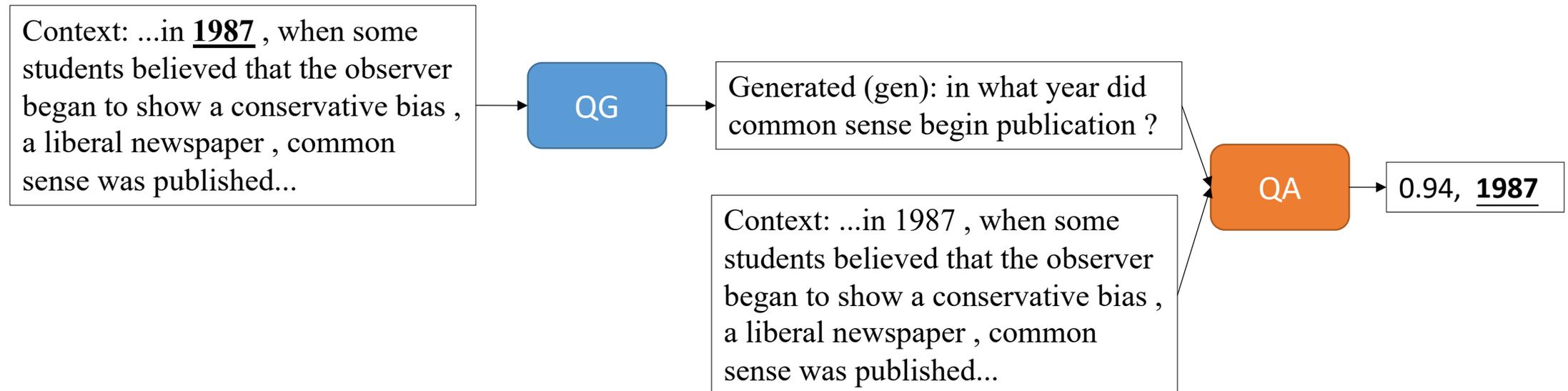


Question Generation with Semantic Validity Knowledge



- QAP (**Q**uestion **A**nswering **P**robability) reward:
 - From QA (**Q**uestion **A**nswering) model
 - Represents “the probability that the generated question can be correctly answered by the given answer”

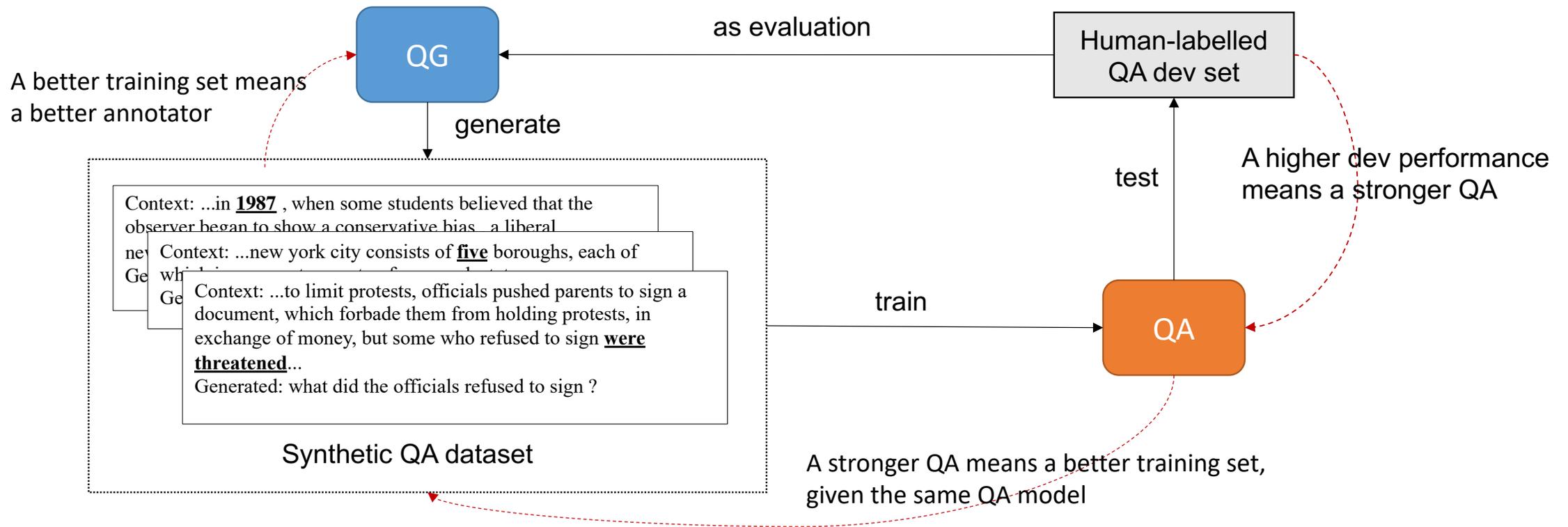
$$p_{qa}(a|q_{gen}, context); q_{gen} \sim p_{qg}(q|a, context)$$



Evaluation for QG



- QA-based QG evaluation: Measure the QG model's ability to mimic human annotators in generating QA training data.

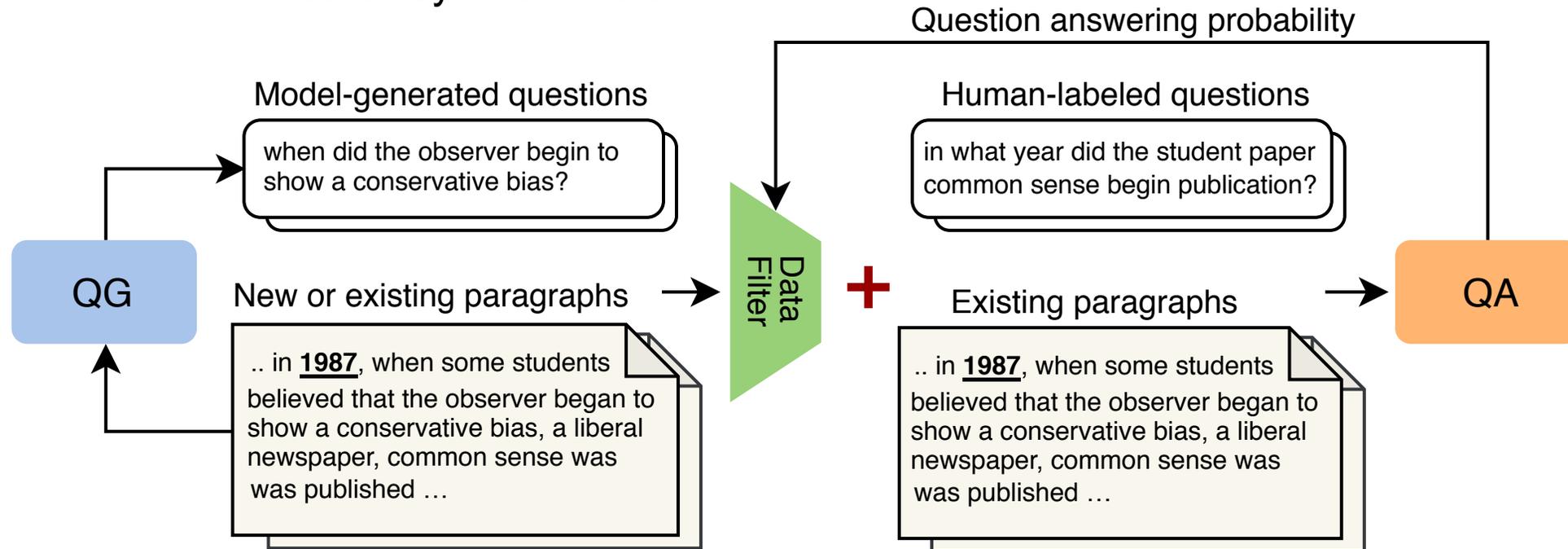


Semi-supervised QA



Augment QA dataset with QG-generated examples (Generate from Existing Articles, and Generate from New Articles)

- (1) QAP filter: To filter out poorly-generated examples; Filter synthetic examples with $QAP < \epsilon$.
- (2) Mixing mini-batch training: To make sure that the gradients from ground-truth data are not overwhelmed by synthetic data, for each mini-batch, we combine half mini-batch ground-truth data with half mini-batch synthetic data.

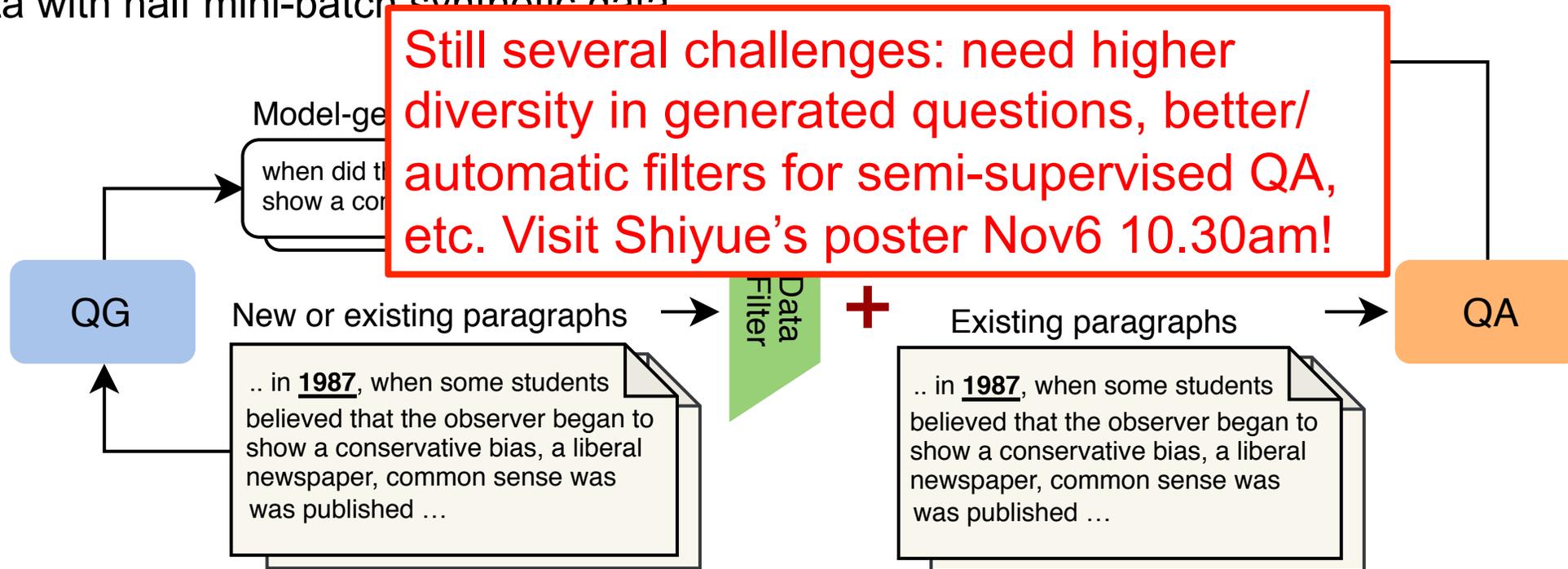


Semi-supervised QA



Augment QA dataset with QG-generated examples (Generate from Existing Articles, and Generate from New Articles)

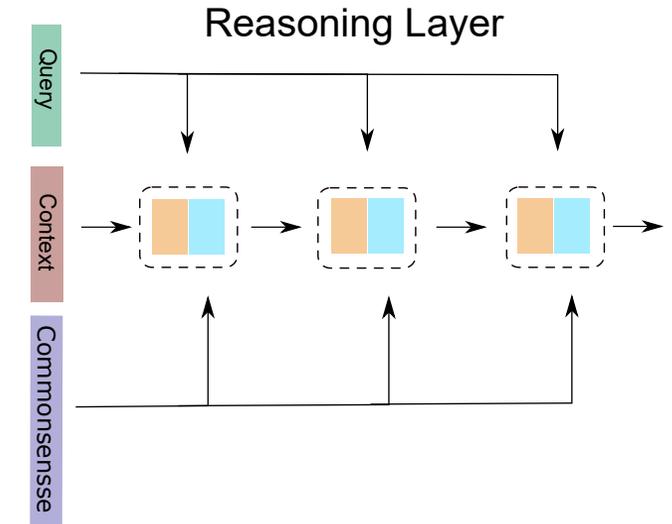
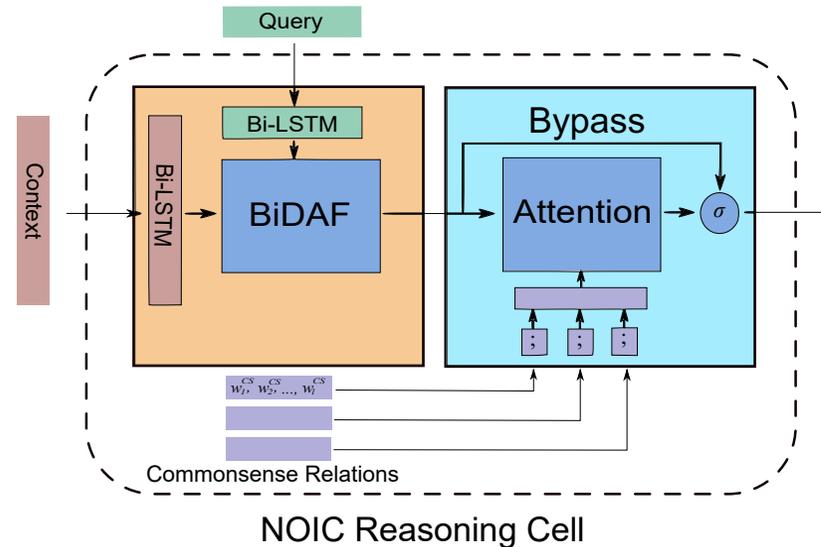
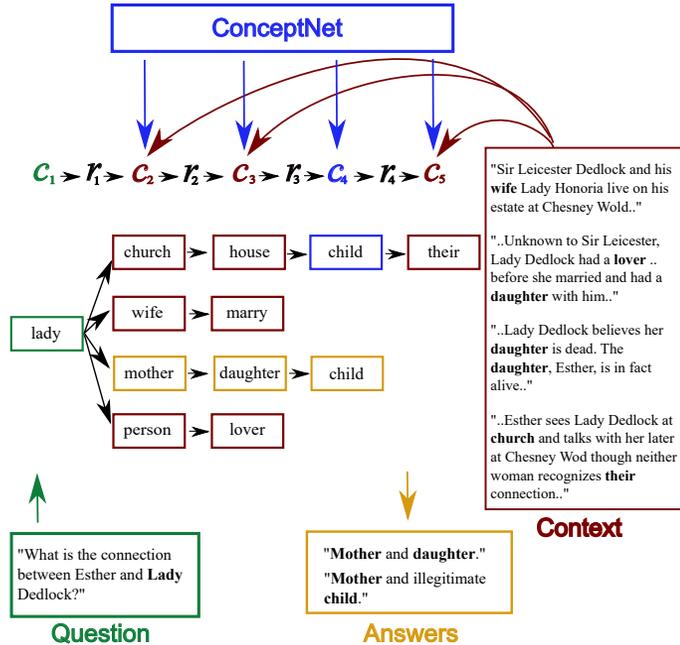
- (1) QAP filter: To filter out poorly-generated examples; Filter synthetic examples with $QAP < \epsilon$.
- (2) Mixing mini-batch training: To make sure that the gradients from ground-truth data are not overwhelmed by synthetic data, for each mini-batch, we combine half mini-batch ground-truth data with half mini-batch synthetic data



Commonsense in Generative Q&A Reasoning



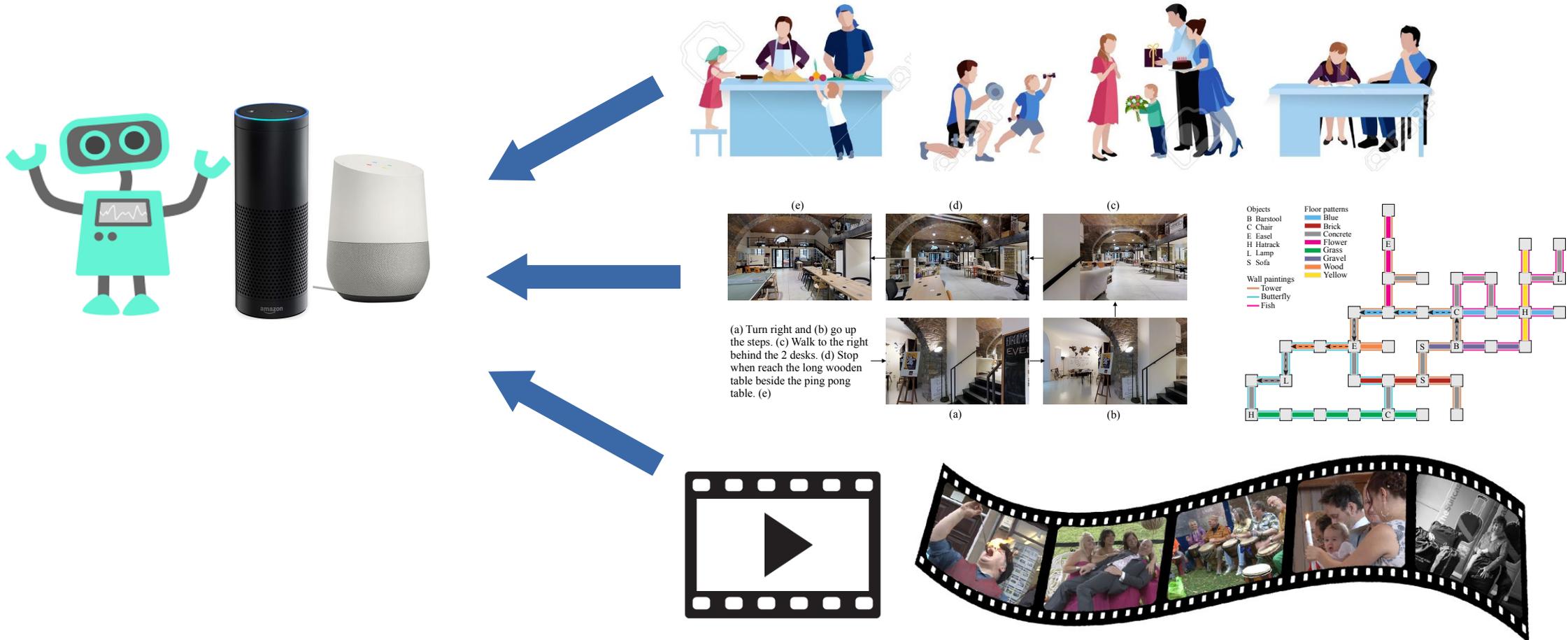
- We use 'bypass-attention' mechanism to reason jointly on both internal context and external commonsense, and essentially learn when to fill 'gaps' of reasoning and with what information



Part2: Spatial, Video-Grounded NLG/Dialogue Models



- NLG/dialogue model should “see” daily activities around it and condition on that context for generation; and execute+generate instructions for navigation and assembling/arrangement tasks, for joint human-robot collaboration/task-solving.



Navigational Instruction Generation



Fig. 4. Participants' field of view in the virtual world used for the human navigation experiments.

Input: map and path

Floor patterns:

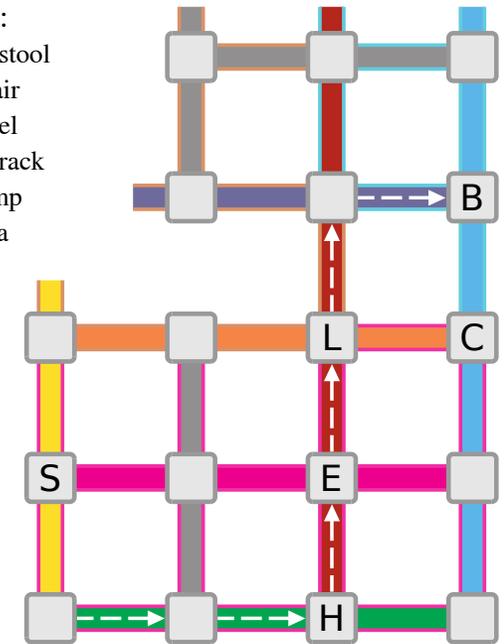
- Blue
- Brick
- Concrete
- Flower
- Grass
- Black
- Wood
- Yellow

Objects:

- B Barstool
- C Chair
- E Easel
- H Hatrack
- L Lamp
- S Sofa

Wall paintings:

- Tower
- Butterfly
- Fish



Output: route instruction

“turn to face the grass hallway. walk forward twice. face the easel. move until you see black floor to your right. face the stool. move to the stool”

Fig. 1. An example route instruction that our framework generates for the shown map and path.

Navigational Instruction Generation

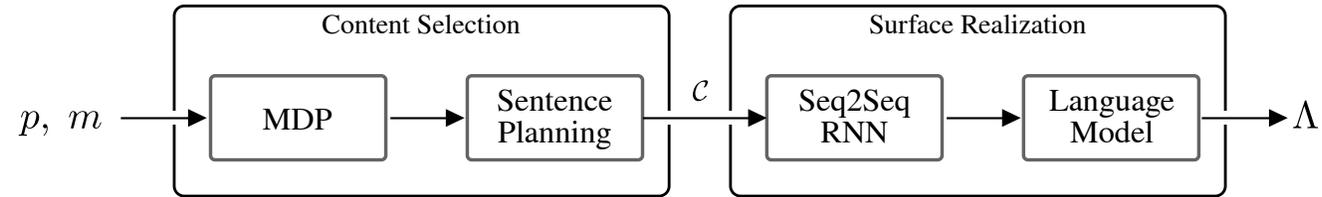


Fig. 2. Our method generates natural language instructions for a given map and path.

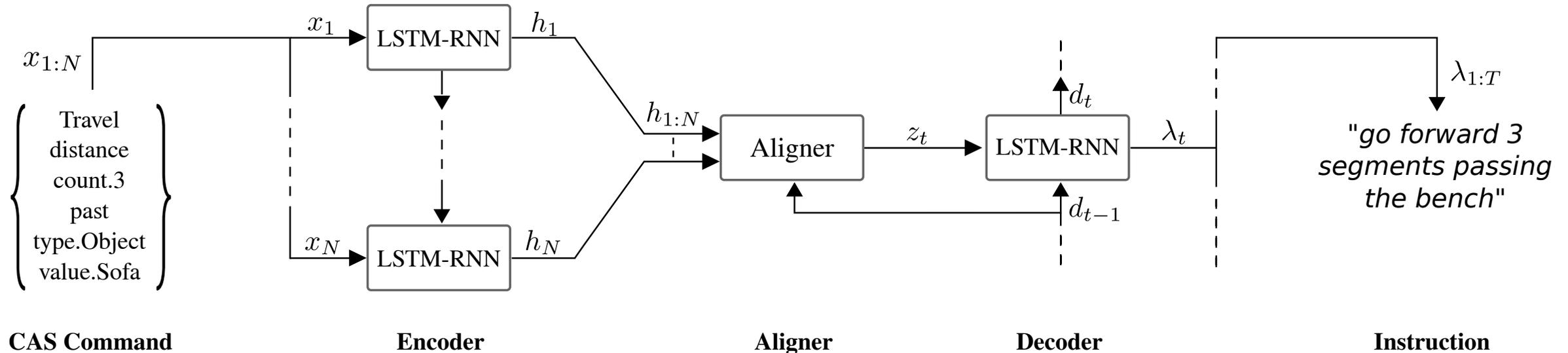


Fig. 3. Our encoder-aligner-decoder model for surface realization.

Navigation Instruction Generation

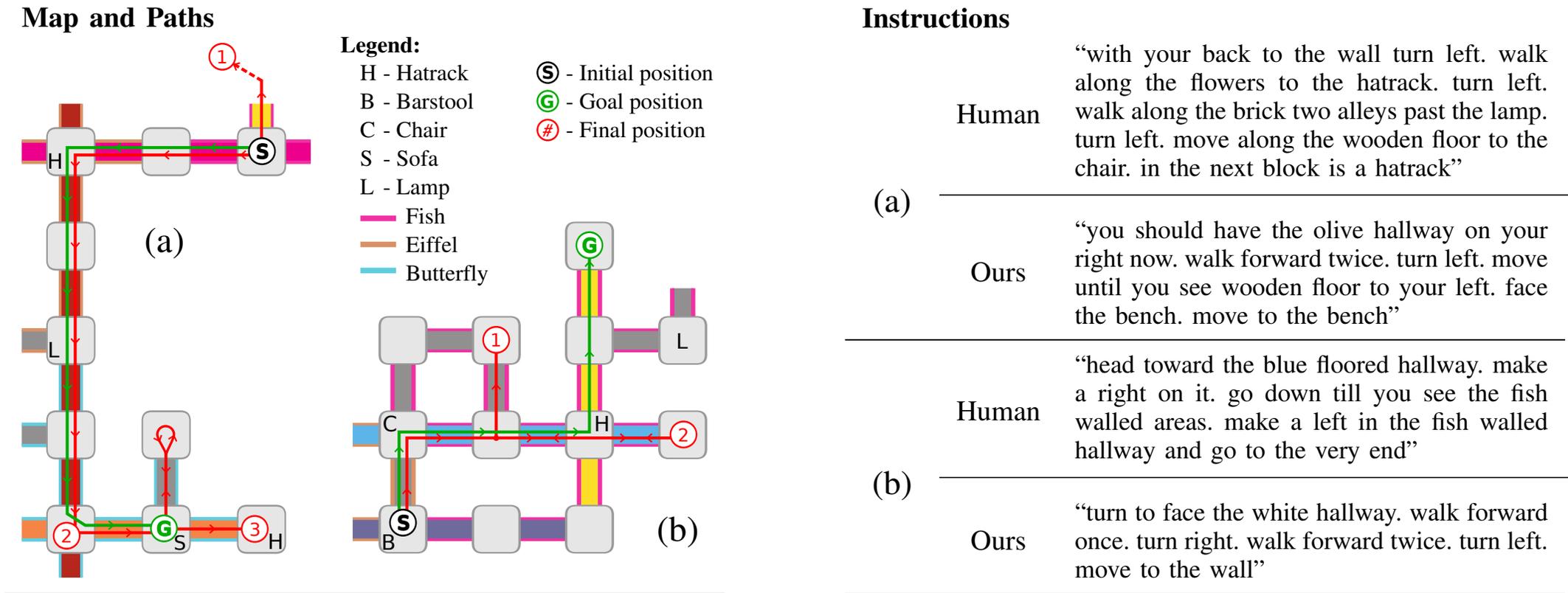


Fig. 8. Examples of paths from the SAIL corpus that ten participants (five for each map) followed according to instructions generated by humans and by our method. Paths in red are those traversed according to human-generated instructions, while paths in green were executed according to our instructions. Circles with an “S” and “G” denote the start and goal locations, respectively.

Room-to-Room Navigation with Instruction Generation



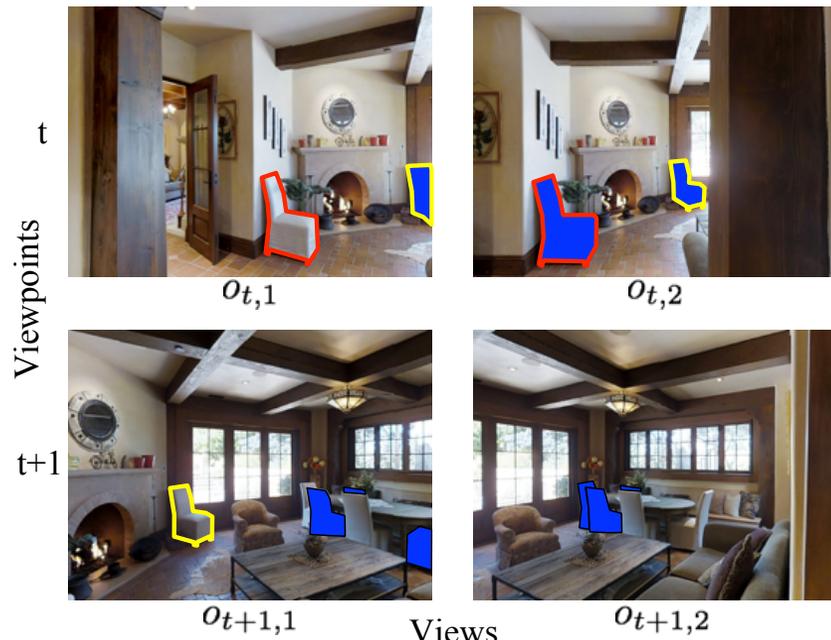
- Learning to Navigate Unseen Environments: Back Translation with Environmental Dropout (to create new rooms with view and viewpoint consistency; generate instructions for new rooms; use generated room-instruction data in semi-supervised setup)



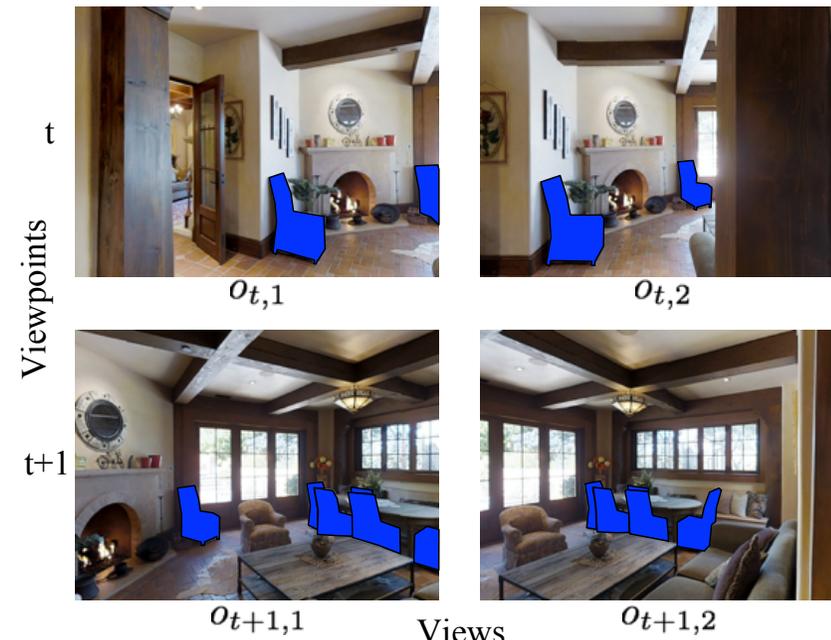
Room-to-Room Navigation with Instruction Generation



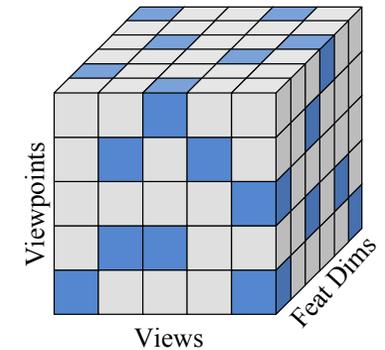
- Learning to Navigate Unseen Environments: Back Translation with Environmental Dropout (to create new rooms with view and viewpoint consistency; generate instructions for new rooms; use generated room-instruction data in semi-supervised setup)



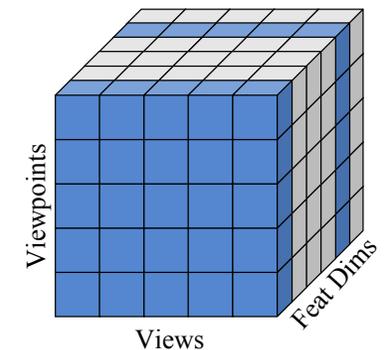
(a) Feature dropout



(b) Environmental dropout



(a) Feature dropout

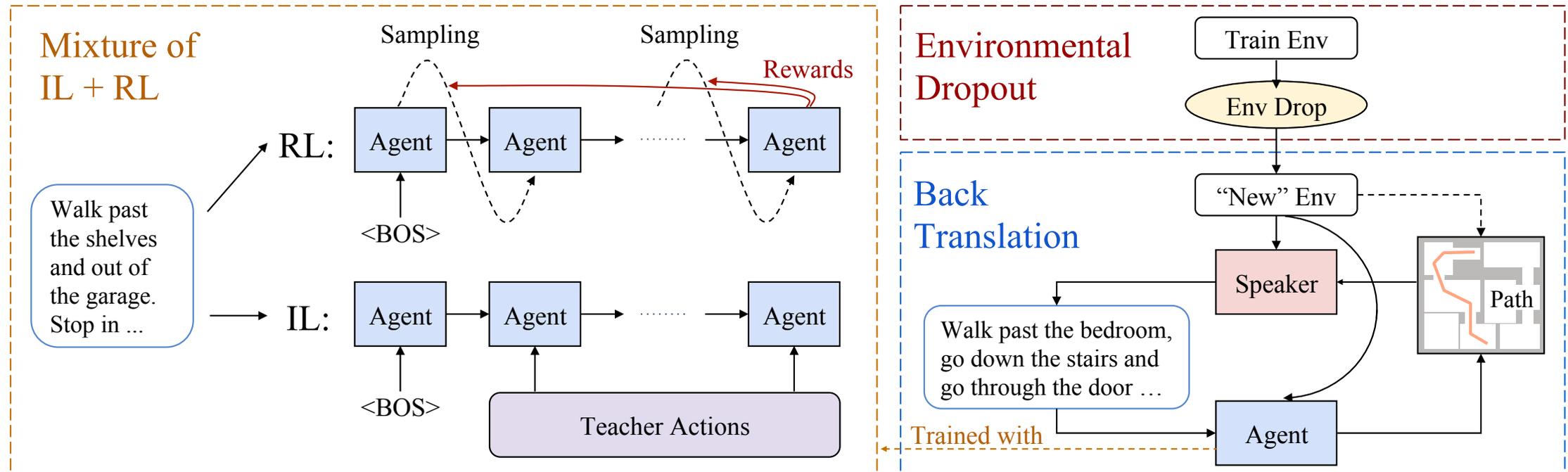


(b) Environmental dropout

Room-to-Room Navigation with Instruction Generation



- Learning to Navigate Unseen Environments: Back Translation with Environmental Dropout (to create new rooms with view and viewpoint consistency; generate instructions for new rooms; use generated room-instruction data in semi-supervised setup)



Room-to-Room Navigation with Instruction Generation



EvalAI - All Challenges Forum Sign Up [Log In](#)

Baseline submission

Rank	Participant team	length	error	oracle success	success	spl	Last submission at
1	human	11.85	1.61	0.90	0.86	0.76	1 year ago
2	Back Translation with Environmental Dropout (with Beam Search) (null)	686.82	3.26	0.99	0.69	0.01	10 months ago
3	vBot (Greedy)	10.24	3.76	0.71	0.65	0.62	3 months ago
4	Back Translation with Environmental Dropout (exploring unseen environments before testing)	9.79	3.97	0.70	0.64	0.61	10 months ago
5	Reinforced Cross-Modal Matching (optimized for SR; with beam search)	357.62	4.03	0.96	0.63	0.02	10 months ago
6	sjtu_test (null)	1,228.45	3.98	0.97	0.62	0.01	10 months ago
7	Self-Monitoring Navigation Agent (with beam search) (Self-Aware Co-Grounded Model)	373.09	4.48	0.97	0.61	0.02	1 year ago
8	Tactical Rewind - long	196.53	4.29	0.90	0.61	0.03	9 months ago
9	Reinforced Cross-Modal Matching + SIL (exploring unseen environments before testing) (SIL-R2)	9.48	4.21	0.67	0.60	0.59	10 months ago
10	AAEI-Agent	13.16	4.61	0.65	0.57	0.50	2 months ago
11	test-sf	10.99	4.57	0.65	0.57	0.50	5 months ago
12	PreSS (Greedy)	10.52	4.53	0.63	0.57	0.53	4 months ago
13	tourist (null)	1,214.94	4.57	0.96	0.56	0.01	11 months ago
14	Tactical Rewind - short	22.08	5.14	0.64	0.54	0.41	10 months ago
15	Speaker-Follower (optimized for success rate) (Speaker-Follower)	1,257.38	4.87	0.96	0.53	0.01	1 year ago
16	Kjtest-sp	948.16	4.91	0.95	0.53	0.01	7 months ago
17	licr19	13.05	5.14	0.60	0.51	0.45	4 months ago

Room-to-Room Navigation with Instruction Generation

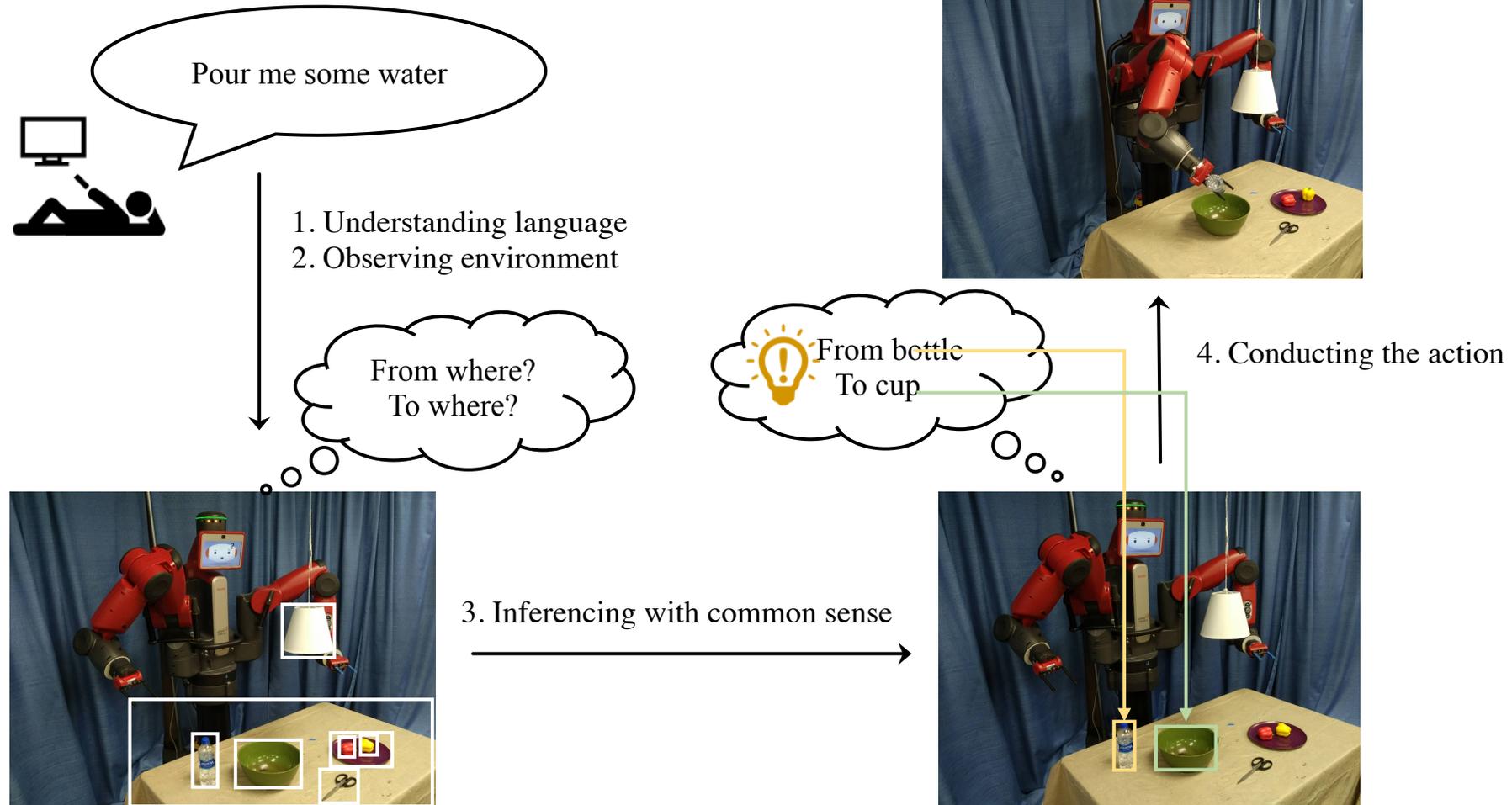


EvalAI - All Challenges Forum Sign Up Log In

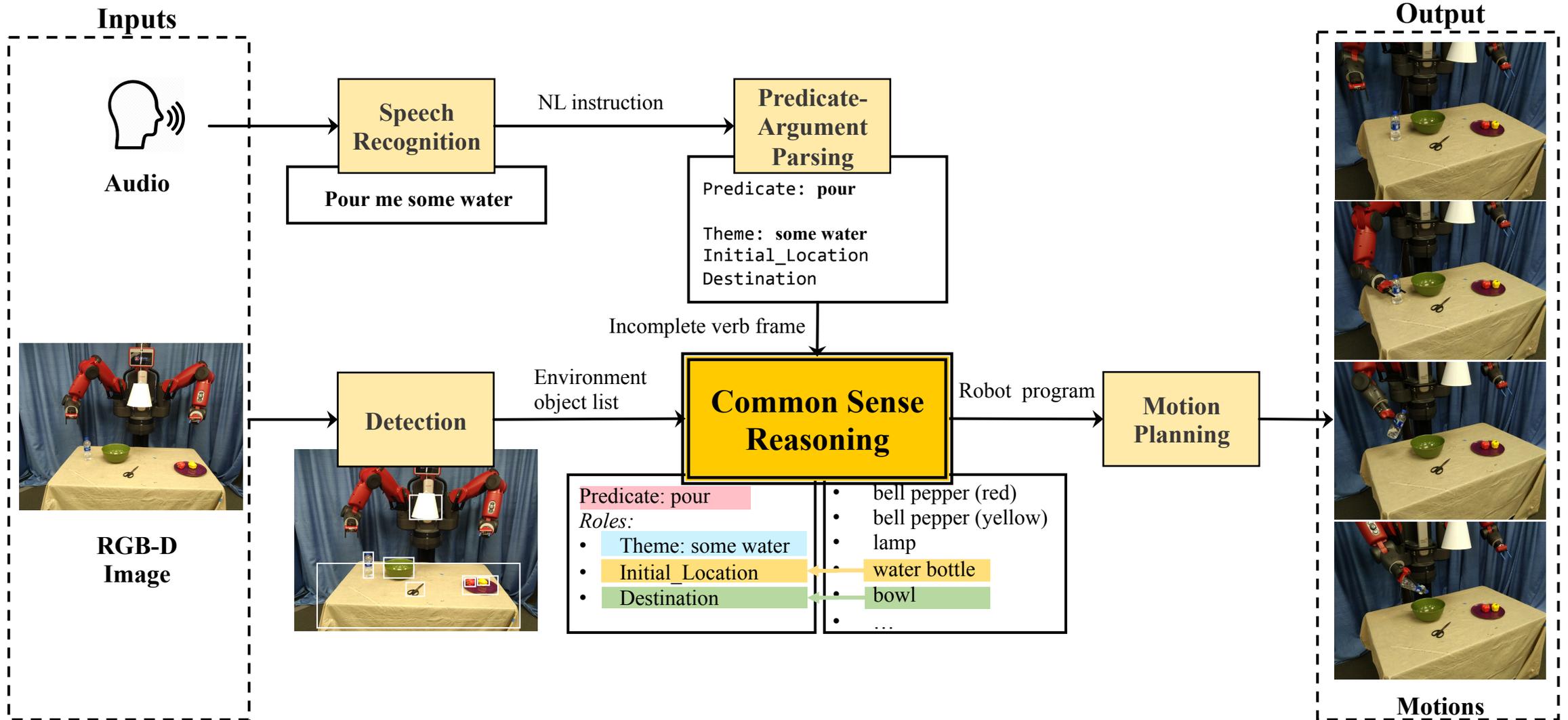
Rank	Participant team	length	error	oracle success	success	spl	Last submission at
1	human	11.85	1.61	0.90	0.86	0.76	1 year ago
2	Back Translation with Environmental Dropout (with Beam Search) (null)	686.82	3.26	0.99	0.69	0.01	10 months ago
3	vBot (Greedy)	10.24	3.76	0.71	0.65	0.62	3 months ago
4	Back Translation with Environmental Dropout (exploring unseen environments before testing)	9.79	3.97	0.70	0.64	0.61	10 months ago
5	Reinforced Cross-Modal Matching (optimized for SR; with beam search)	357.62	4.03	0.96	0.63	0.02	10 months ago
6	sjtu_test (null)	1,228.45	3.98	0.97	0.62	0.01	10 months ago
7	Self-Monitoring Navigation Agent (with beam search) (Self-Aware Co-Grounded Model)	373.09	4.48	0.97	0.61	0.02	1 year ago
8	Tactical Rewind - long	196.53	4.29	0.90	0.61	0.03	9 months ago
9	Reinforced Cross-Modal Matching + SIL (exploring unseen environments before testing) (SIL-R2)	9.48	4.21	0.67	0.60	0.59	10 months ago
10	AAEI-Agent	13.16	4.61	0.65	0.57	0.50	2 months ago
11	test-sf	10.99	4.57	0.65	0.57	0.50	5 months ago
12	PreSS (Greedy)	10.52	4.53	0.63	0.57	0.53	4 months ago
13	tourist (null)	1,214.94	4.57	0.96	0.56	0.01	11 months ago
14	Tactical Rewind - short	22.00	5.14	0.64	0.54	0.41	10 months ago
15	Speaker-Follower (optimized for success rate) (Speaker-Follower)	10.00	4.50	0.60	0.50	0.50	10 months ago
16	Kjtest-sp	10.00	4.50	0.60	0.50	0.50	10 months ago
17	licr19	10.00	4.50	0.60	0.50	0.50	10 months ago

Still several challenges/ long way to go, e.g., better object detectors, diverse language, etc.!

Commonsense via Robotic Instruction Completion



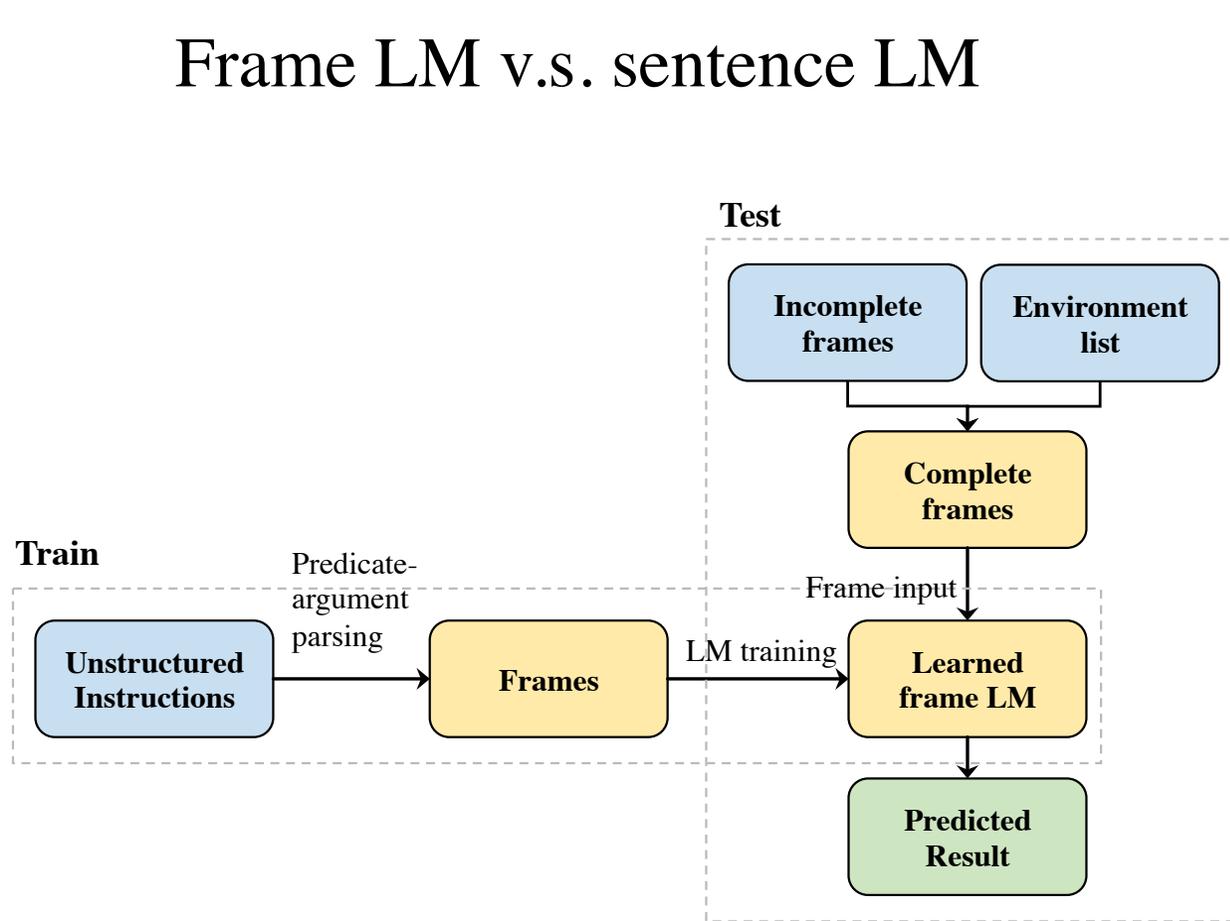
Commonsense via Robotic Instruction Completion



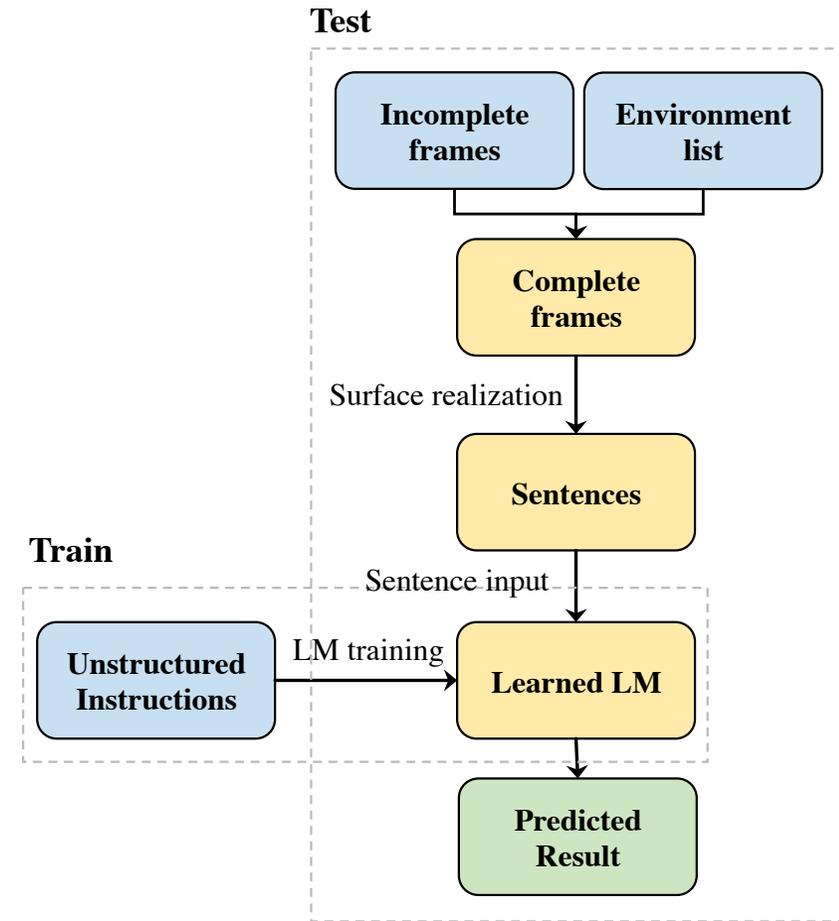
Commonsense via Robotic Instruction Completion



Frame LM v.s. sentence LM



Frame LM



Sentence LM

Commonsense via Robotic Instruction Completion

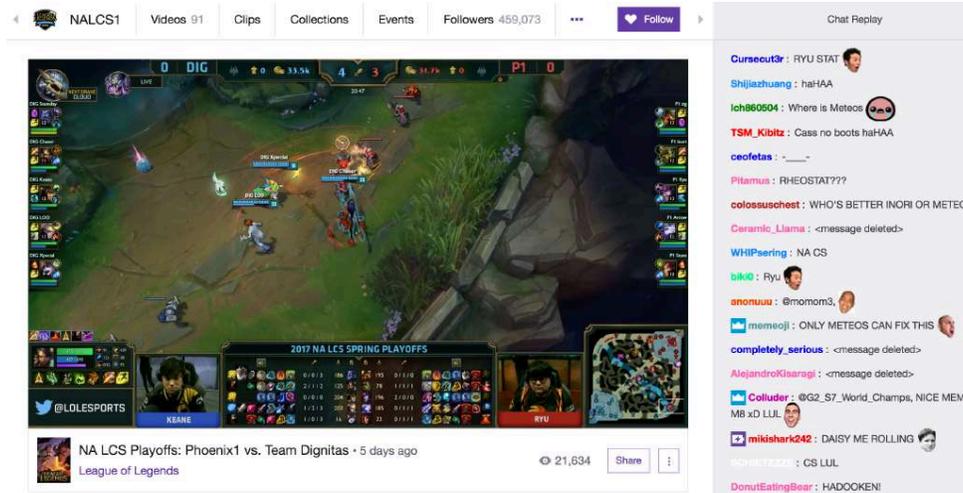


<https://drive.google.com/file/d/1C9xsuyW1bVBzLimvVFbBfOcKCzV5ueHs/view>

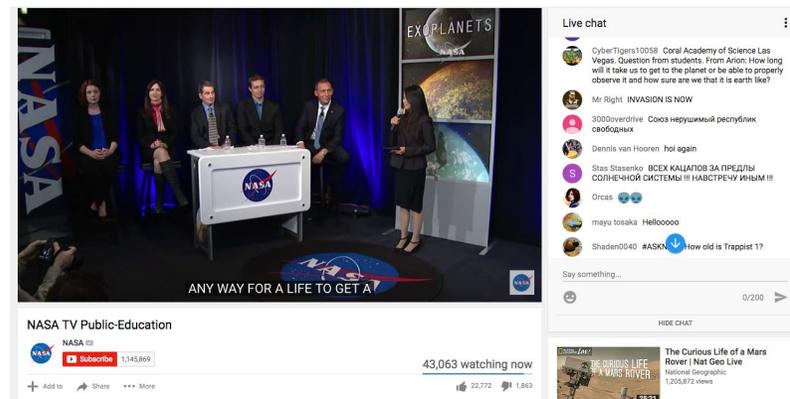
New Spatio-Temporal Video+Dialogue Task



- Video + Chat: conversations grounded in concrete video events!



(a) Twitch



(b) Youtube

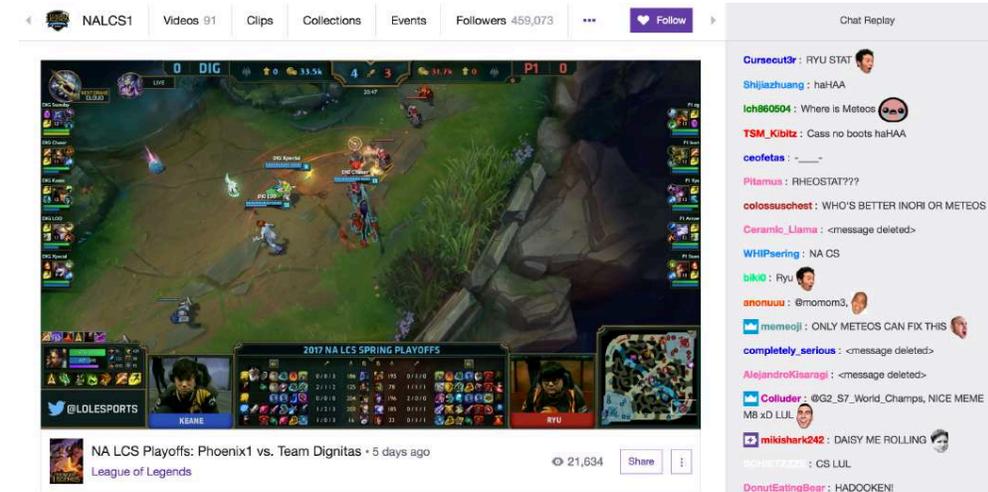


(c) Facebook

New Spatio-Temporal Video+Dialogue Task



- Very interesting chat language!
 - Time-constrained, not just space
 - Lots of special vocab, symbols, emoticons
 - Multi-user with several interleaving turns
 - Multi-lingual

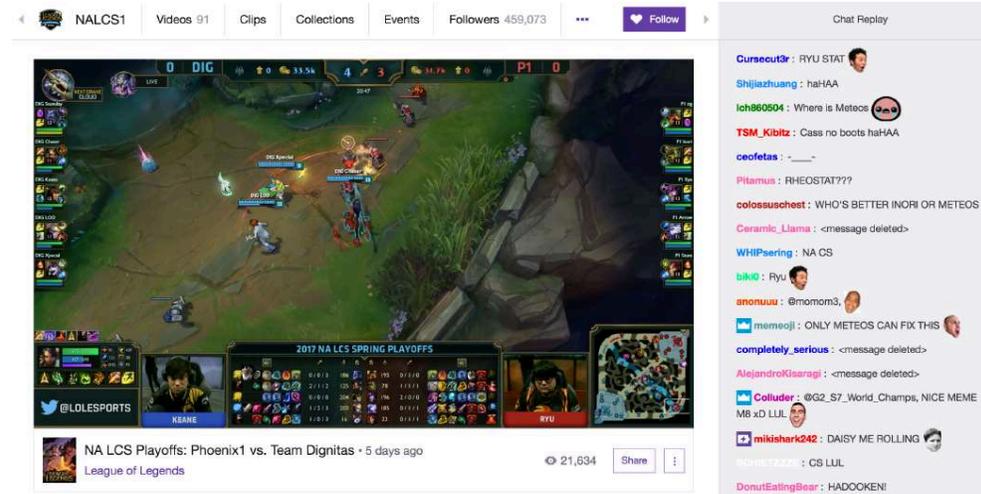


Code/Data: <https://github.com/chengyangfu/Pytorch-Twitch-LOL>

New Spatio-Temporal Video+Dialogue Task

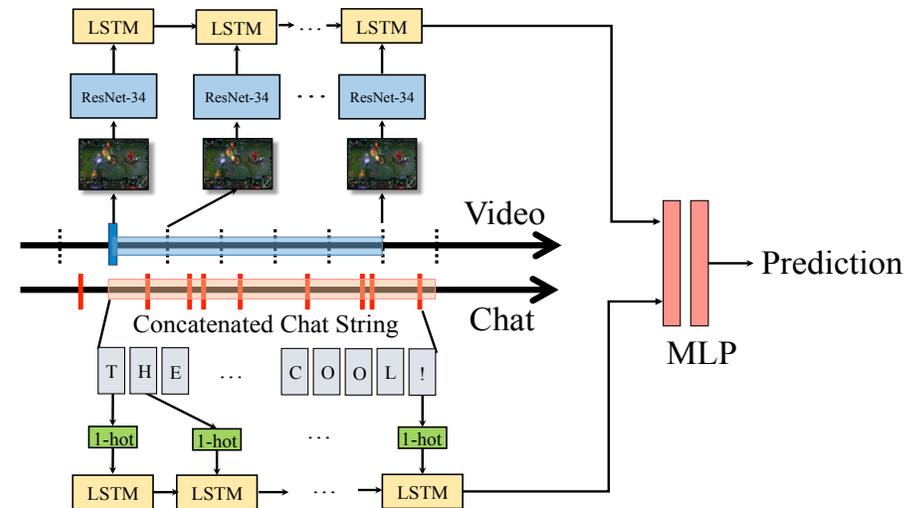


- Very interesting chat language!
 - Time-constrained, not just space
 - Lots of special vocab, symbols, emoticons
 - Multi-user with several interleaving turns
 - Multi-lingual
- First, we predicted the summary/highlight frames of the full video using joint features from video and user reactions from chat dialogue in English +Chinese (via character-level model to capture the new language style/formats)



Method	Data	NALCS	LMS
L-Char-LSTM	chat	43.2	39.7
V-CNN-LSTM	video	72.2	69.2
<i>lv</i> -LSTM	chat+video	74.7	70.0

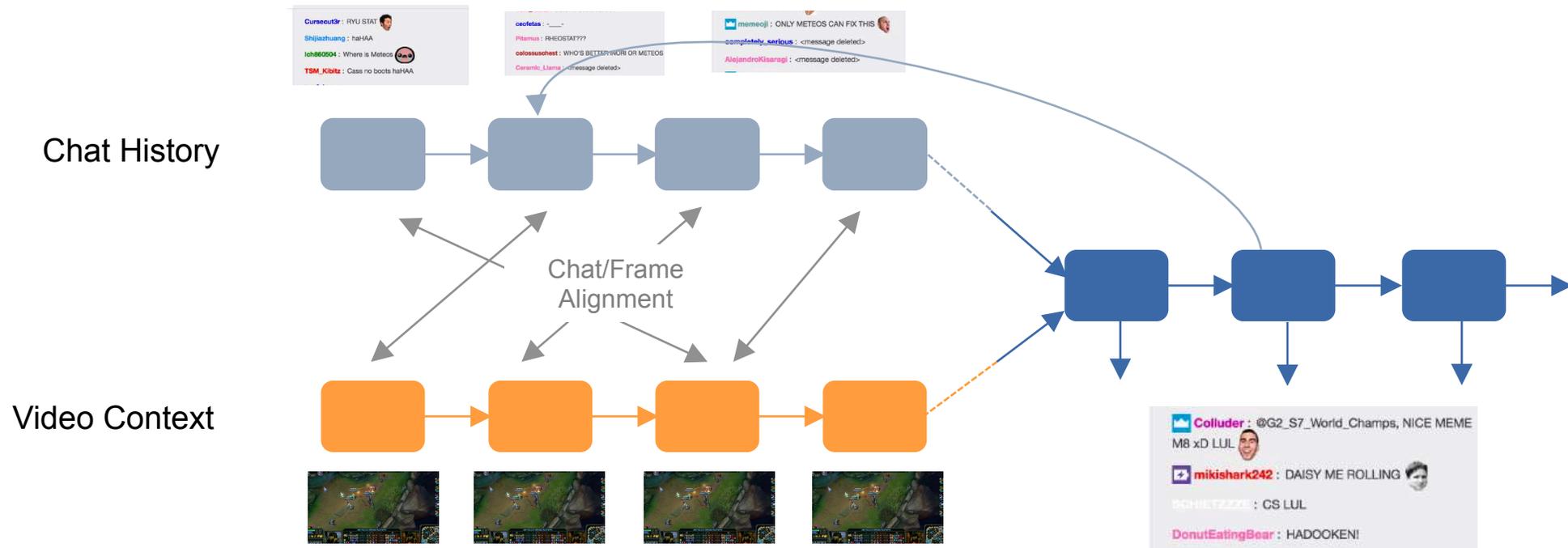
Table 3: Test Results on the NALCS (English) and LMS (Traditional Chinese) datasets.



Dialogue Generation on Video Context



- Next: Generating chat responses given the video and previous dialogue history!



Code/Data: <https://github.com/ramakanth-pasunuru/video-dialogue>

Dialogue on Video Context



S1: what an offside trap
OMEGALUL

S2: Lol that finish bro

S3: suprised you didn't
do the extra pass

S4: @S10 a drunk bet?

S5: @S11 thanks mate

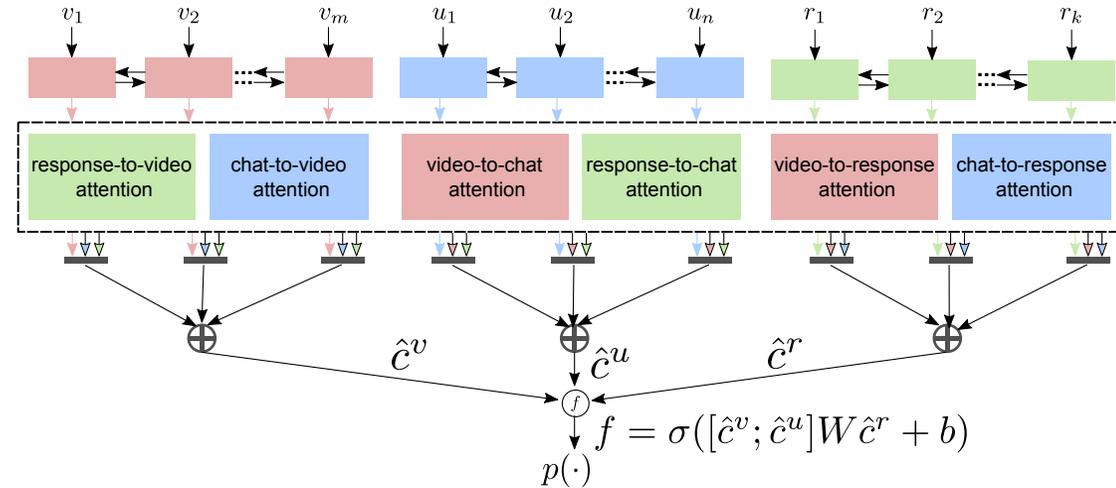
S6: could have passed
one more

S7: Pass that

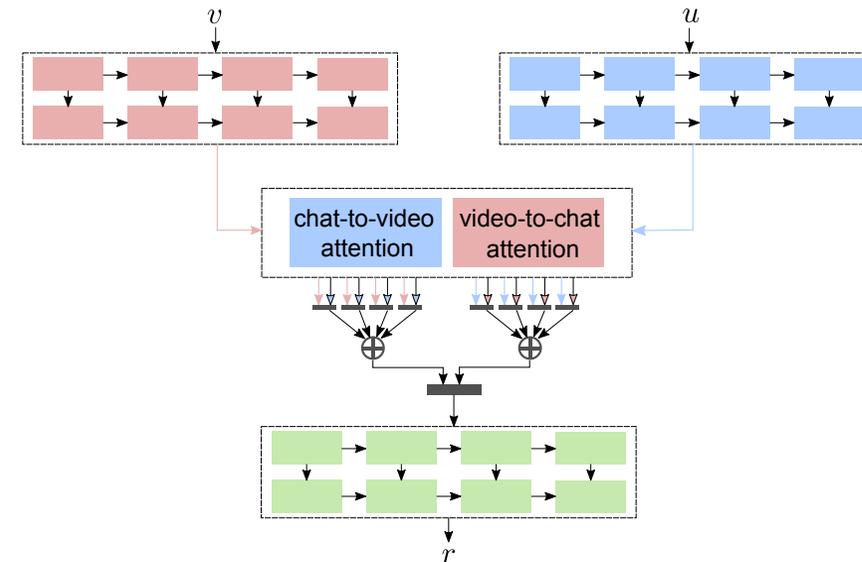
S1: record now!

S8: !record

S9: done a nother pass there



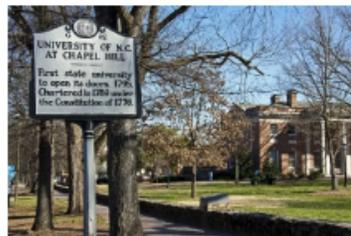
Models	r@1	r@2	r@5
BASELINES			
Most-Frequent-Response	10.0	16.0	20.9
Naive Bayes	9.6	20.9	51.5
Logistic Regression	10.8	21.8	52.5
Nearest Neighbor	11.4	22.6	53.2
Chat-Response-Cosine	11.4	22.0	53.2
DISCRIMINATIVE MODEL			
Dual Encoder (C)	17.1	30.3	61.9
Dual Encoder (V)	16.3	30.5	61.1
Triple Encoder (C+V)	18.1	33.6	68.5
TriDAF+Self Attn (C+V)	20.7	35.3	69.4
GENERATIVE MODEL			
Seq2seq +Attn (C)	14.8	27.3	56.6
Seq2seq +Attn (V)	14.8	27.2	56.7
Seq2seq + Attn (C+V)	15.7	28.0	57.0
Seq2seq + Attn + BiDAF (C+V)	16.5	28.5	57.7



Thoughts/Challenges/Future Work



- Other axes of NLG:
 - Personality (we have done some work on politeness/rudeness- and humor-based language generation)
 - Speed and scalability (hybrid extractive+abstractive summarization with RL connector; SotA +20x speedup)
- Extending the video-dialogue and video-QA models to multiple other languages
- AutoAugment design for other NLG tasks
- More structured commonsense for other NLG tasks
- Better AutoAugment algorithms for speed, input-awareness, RL instability and reward sparsity
- Richer spatial world benchmarks with instruction generation/dialogue



Welcome to the UNC-NLP Research Group

Our lab has research interests in statistical natural language processing and machine learning, with a focus on multimodal, grounded, and embodied semantics (i.e., language with vision and speech, for robotics), human-like language generation and Q&A/dialogue, and interpretable and structured deep learning. We are a group of PhD, MS, BS, and visiting students who work with [Prof. Mohit Bansal](#) and collaborators in the [Computer Science department](#) (lab located in [Brooks Building FB-241C](#)) at the [University of North Carolina \(UNC\) Chapel Hill](#).

News

Aug 2019 Congrats to [Peter Hase](#) for the [Royster Society PhD Fellowship!](#)

Aug 2019 5 new [papers](#) in [EMNLP 2019](#).

July 2019 Congrats to Hyounghun for [ACL 2019 Best Short Paper Nomination!](#)

July 2019 We have a [Postdoc opening](#) - please apply!

July 2019 Thanks for the [NSF-CAREER Award \(details\)](#).

July 2019 Thanks for the [Google Focused Research Award \(details\)](#).

May 2019 6 new [papers](#) in [ACL 2019](#).

Apr 2019 Congrats to [Darryl Hannan](#) for the 3-year [NSF PhD Fellowship!](#)

Mar 2019 Congrats to [Hao Tan](#) for [1st Rank](#) on the Room-to-Room Vision-Language-Navigation Leaderboard!

Feb 2019 5 new [papers](#): 3 in [NAACL 2019](#), 1 in [CVPR 2019](#), 1 in [ICRA 2019](#).

Jan 2019. Congrats to [Ramakanth Pasunuru](#) for being awarded the 2-year [Microsoft Research PhD Fellowship!](#)

Mar 2018. Thanks to Adobe for the [Adobe Research Award](#).

Feb 2018. [9 new 2018 papers](#) in NAACL, CVPR, AAAI, WACV.

Sept 2017. Thanks to DARPA for the [DARPA Young Faculty Award \(link\)](#).

Sept 2017. Thanks to Facebook for the [Facebook ParLAI Research Award](#).

July 2017. 3 papers at [EMNLP 2017](#) and 2 papers at the [Summarization-Frontiers](#) and [RepEval](#) workshops.

June 2017. Top single model results on the [RepEval-NLI Shared Task](#) at EMNLP 2017 (congrats Yixin!).

June 2017. [Outstanding Paper Award](#) at ACL 2017 (congrats Ram!).

Feb 2017. Thanks to Google for a [Google Faculty Research Award \(link\)](#).

Nov 2016. 3 papers on [navigational instruction generation, coherent dialogue w/ attn-LMs](#), and on [context-RNN-GAN models](#) to appear at [AAAI 2017](#) and [HRI 2017](#).

July 2016. [5 papers](#) to appear at [EMNLP 2016](#): visual story sorting, visual question relevance, neural network interpretation (for

Tweets by @uncnlp

UNC NLP Retweeted



emnlp2019

@emnlp2019

Registration for EMNLP 2019 will open in a few days. In the meantime, you can have a look at the registration fees for the conference. [emnlp-ijcnlp2019.org/registration/](#)

EMNLP-IJCNLP 2019 Registration Fees

Type	Register	Full package	Main conference	Main + 1 Day	1 Day	2 Days
Regular	Early	\$995	\$685	\$825	\$220	\$330
	Late	\$1120	\$800	\$1075	\$275	\$415
	Onsite	\$1315	\$905	\$1235	\$330	\$495

PhD Students



Lisa Bauer
PhD at UNC



Darryl Hannan
PhD at UNC



Peter Hase
PhD at UNC



Yichen Jiang
PhD at UNC



Hyounghun Kim
PhD at UNC
(co-advised w/ H. Fuchs)



Jie Lei
PhD at UNC
(co-advised w/ T. Berg)



Adyasha Maharana
PhD at UNC



Yixin Nie
PhD at UNC



Ramakanth Pasunuru
PhD at UNC



Swarnadeep Saha
PhD at UNC



Hao Tan
PhD at UNC



Shiyue Zhang
PhD at UNC



Yubo Zhang
PhD at UNC
(co-advised w/ A. Tropsha)



Xiang Zhou
PhD at UNC

Undergraduate Students



Tsion Coulter
UG at UNC



Han Guo
UG at UNC



Akshay Jain
UG at UNC



Sweta Karlekar
UG at UNC



Antonio Mendoza
UG at UNC



Yicheng Wang
UG at UNC



Songhe Wang
UG at UNC



Thank you!

Webpage: <http://www.cs.unc.edu/~mbansal/>

Email: mbansal@cs.unc.edu

UNC-NLP Lab: <http://nlp.cs.unc.edu/>

Postdoc Openings!!! [~mbansal/postdoc-advt-unc-nlp.pdf](http://www.cs.unc.edu/~mbansal/postdoc-advt-unc-nlp.pdf)