

Deep Reinforcement Learning-based Image Captioning with Embedding Reward

Zhou Ren¹, Xiaoyu Wang¹, Ning Zhang¹, Xutao Lv¹, Li-Jia Li²

¹Snap Inc.

²Google Inc.

Decision-Making Framework for Image Captioning

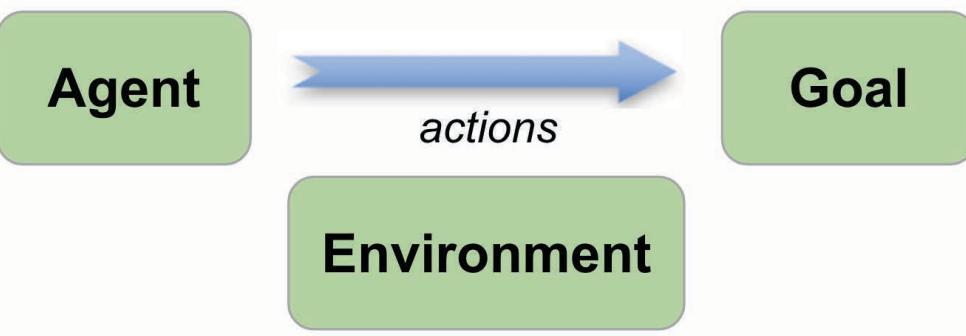
Limitations of the Encoder-Decoder Framework:

- Only **local** information is utilized in training and inference
- Prone to **accumulate** generation errors during inference
- Sensitive** to beam sizes in inference

Our Target for the New Framework:

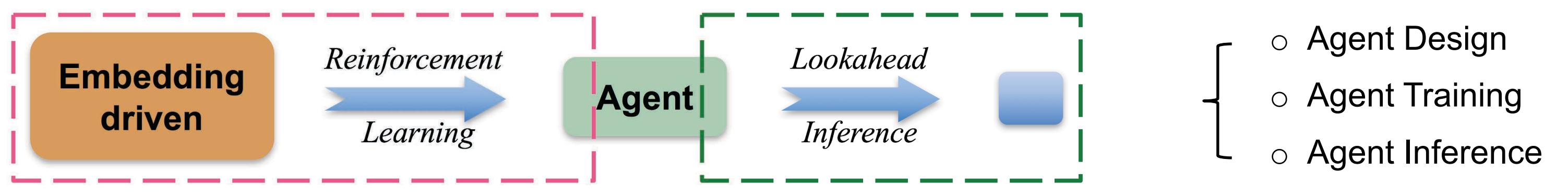
- Better at utilizing **global** information
- Less likely** to accumulate generation errors during inference
- Less sensitive** to beam sizes in inference

Problem Re-formulation in the New Framework:

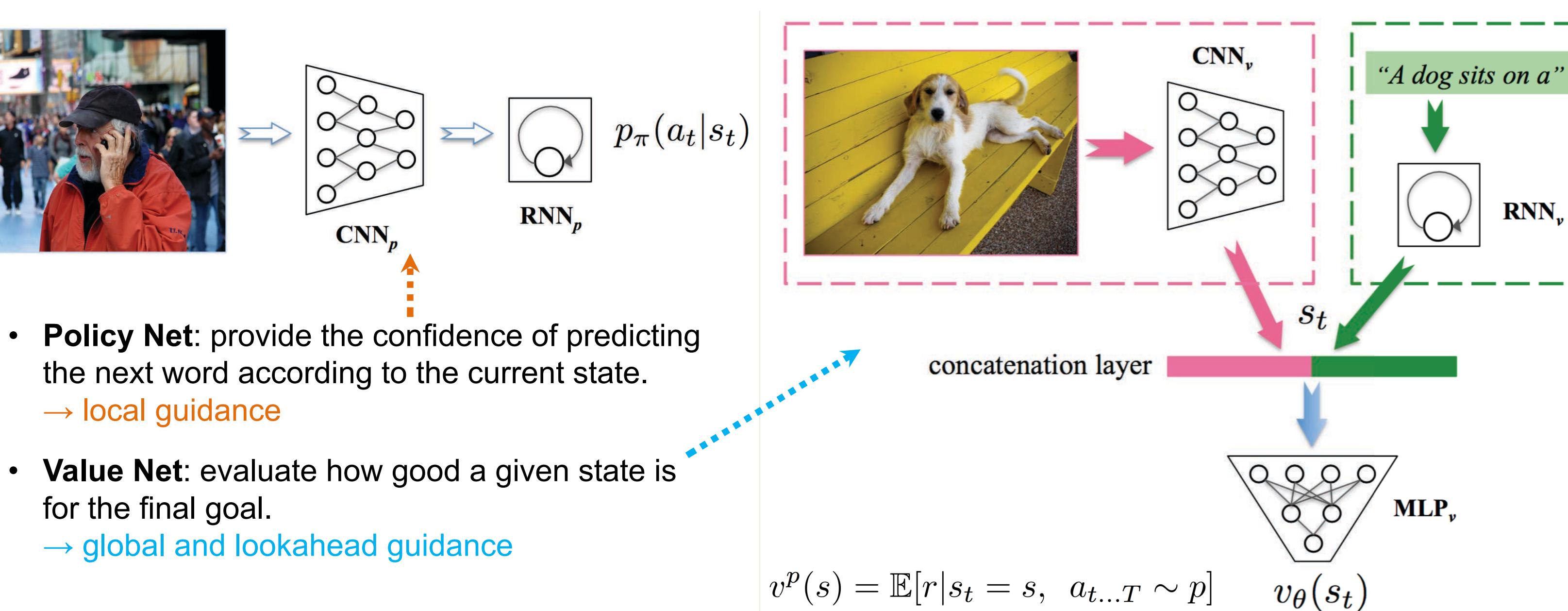


- Agent: the image captioning model to learn
- Environment: an image \mathbf{I} , and the words predicted so far
- State s_t : the representation of environment at each time step t
- Action a_t : the word to predict at each time step t

Overview of Our Approach:

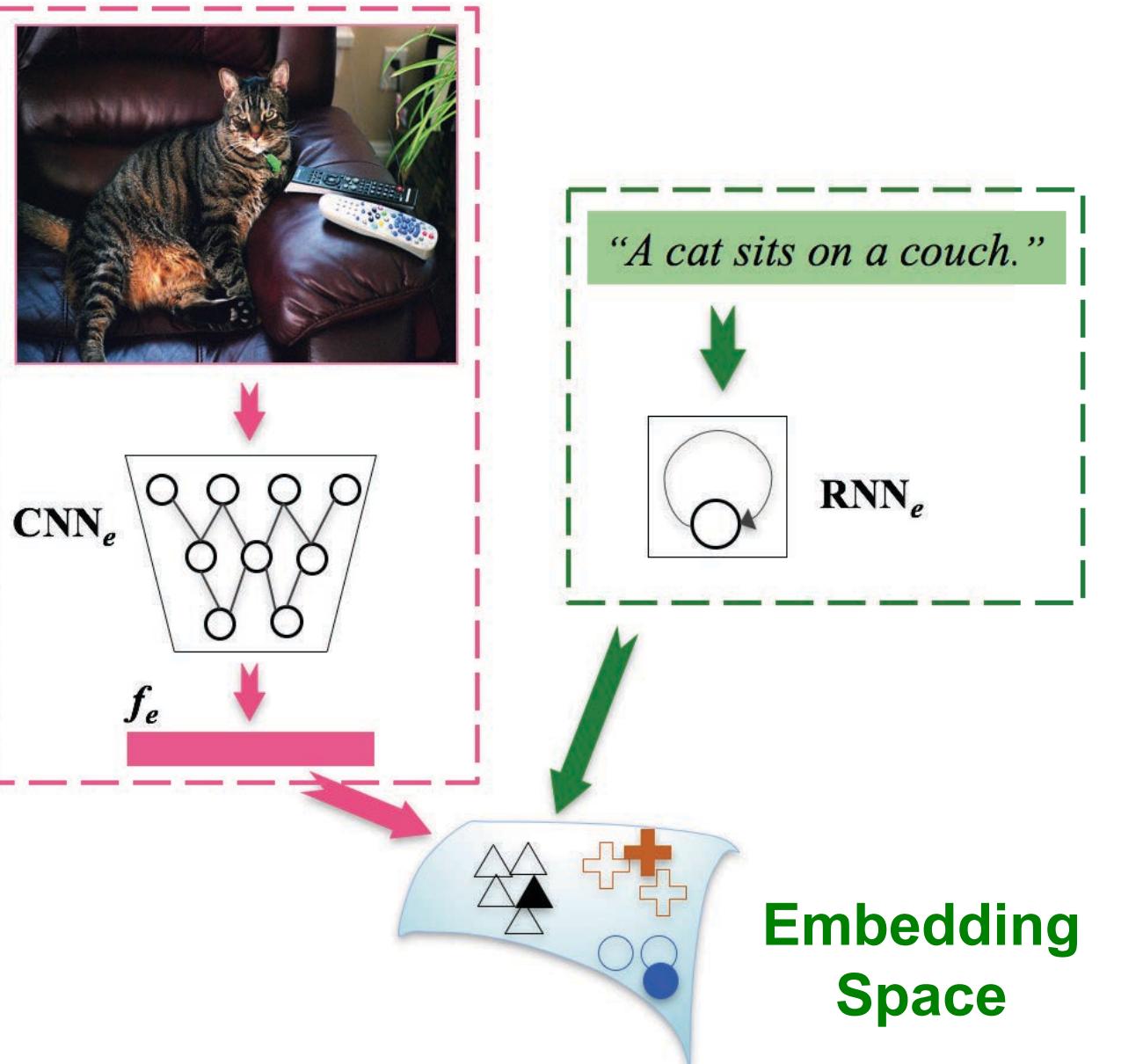


Agent Design: Policy Network + Value Network

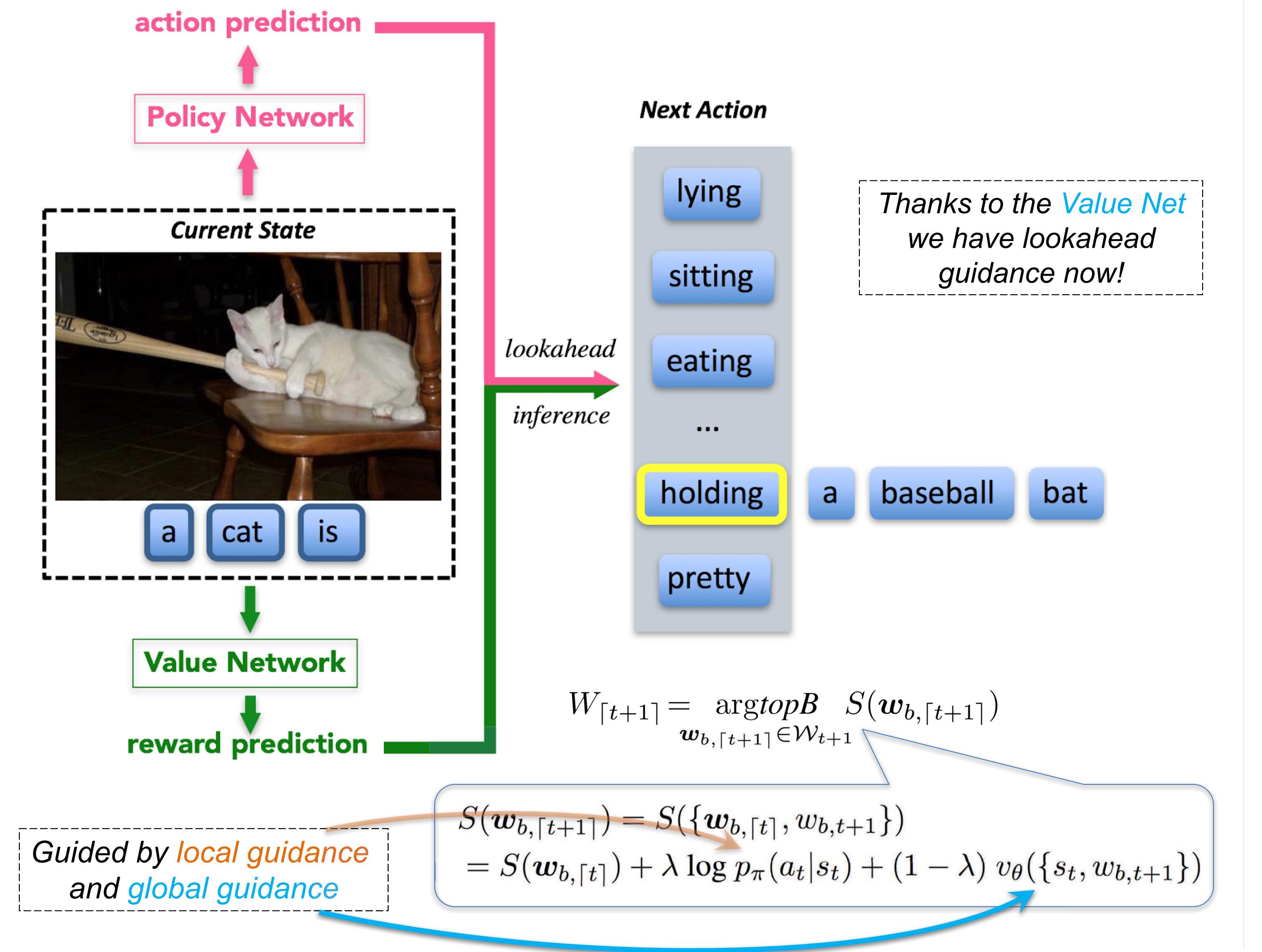


Agent Training: RL with Embedding Reward

- Pretrain Policy Net with cross entropy loss
- Pretrain Value Net with mean squared loss
- Jointly train Policy Net and Value Net using deep Reinforcement Learning
 - an Actor-Critic RL model
 - use MIXER training [Ranzato et al. 2016]
 - reward is defined by visual-semantic embedding



Agent Inference: Lookahead Inference



Experiments on MS-COCO

better at capturing global information



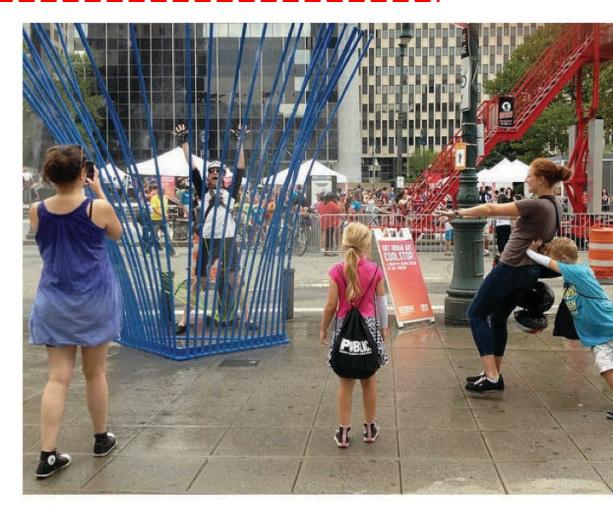
GT: a very large sheep is standing in the grass
SL: a large brown bear walking across a lush green field
Ours: a brown and white sheep standing on a lush green field



GT: the plane is parked at the gate at the airport terminal
SL: a passenger train that is pulling into a station
Ours: a white airplane parked at an airport terminal



GT: a man holding a snowboard next to a man in scary costume
SL: a woman standing in a living room holding a wii controller
Ours: a woman sitting on a ledge holding a snowboard



GT: people are standing outside in a busy city street
SL: a group of young people playing a game of basketball
Ours: a group of people that are standing in the street



GT: a couple of kids walking with umbrellas in their hands
SL: a couple of people that are walking in the snow
Ours: a couple of people walking down a street holding umbrellas



GT: a small dog eating a plate of broccoli
SL: a dog that is sitting on a table
Ours: a dog that is eating some food on a table

Quantitative Results and Ablation Study:

Methods	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	Rouge-L	CIDEr
Google NIC [44]	0.666	0.461	0.329	0.246	—	—	—
m-RNN [30]	0.67	0.49	0.35	0.25	—	—	—
BRNN [17]	0.642	0.451	0.304	0.203	—	—	—
LRCN [7]	0.628	0.442	0.304	0.21	—	—	—
MSR/CMU [3]	—	—	—	0.19	0.204	—	—
Spatial ATT [46]	0.718	0.504	0.357	0.25	0.23	—	—
gLSTM [15]	0.67	0.491	0.358	0.264	0.227	—	0.813
MIXER [35]	—	—	—	0.29	—	—	—
Spatial ATT [48] *	0.709	0.537	0.402	0.304	0.243	—	—
DCC [13] *	0.644	—	—	—	0.21	—	—
Ours	0.713	0.539	0.403	0.304	0.251	0.525	0.937

Validated the contribution of all key components in our framework: embedding; reinforcement learning; and lookahead inference.

Our approach is modular w.r.t. the agent design

Parameter Sensitivity Analysis:

λ	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	Rouge-L	CIDEr
0	0.638	0.471	0.34	0.247	0.233	0.501	0.8
0.1	0.683	0.51	0.373	0.274	0.248	0.516	0.894
0.2	0.701	0.527	0.389	0.288	0.248	0.521	0.922
0.3	0.71	0.535	0.398	0.298	0.251	0.524	0.934
0.4	0.713	0.539	0.403	0.304	0.247	0.525	0.937
0.5	0.71	0.538	0.402	0.304	0.246	0.524	0.934
0.6	0.708	0.535	0.399	0.301	0.245	0.522	0.923
0.7	0.704	0.531	0.395	0.297	0.243	0.52	0.912
0.8	0.704	0.526	0.392	0.295	0.241	0.518	0.903
0.9	0.698	0.524	0.389	0.293	0.24	0.516	0.895
1	0.694	0.52	0.385	0.289	0.238	0.513	0.879

Our approach works well due to both the global and local guidance

Method	Beam size	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	Rouge-L	CIDEr
SL	5	0.696	0.522	0.388	0.29	0.238	0.513	0.876
SL	10	0.692	0.519	0.384	0.289	0.237	0.512	0.872
SL-Embed	7	0.706	0.533	0.395	0.298	0.243	0.52	0.916
SL-RawVN	25	0.683	0.508	0.374	0.271	0.223	0.503	0.850
SL	50	0.680	0.505	0.372	0.279	0.233	0.503	0.849
SL	100	0.679	0.504	0.372	0.279	0.233	0.503	0.849
Ours	5	0.711	0.538	0.403	0.302	0.251	0.524	0.934
Ours	10	0.713	0.539	0.403	0.304	0.251	0.525	0.937
Ours	25	0.709	0.534	0.398	0.299	0.248	0.522	0.928
Ours	50	0.708	0.533	0.397	0.298	0.247	0.52	0.924
Ours	100	0.707	0.531	0.395	0.297	0.244	0.52	0.92

Our approach is less sensitive to beam sizes comparing to encoder-decoder