

Dialogue systems

NOVEMBER 22, 2018

Elena Voita
Yandex Research,
University of Amsterdam
lena-voita@yandex-team.ru

Plan

- Types of dialogue systems
- Task-oriented
- General conversation: generative
- General conversation: retrieval-based
- Hack of the day (real-life story inside!)

An example of a dialogue system

Где я?

Я думаю, примерно тут:
Шотландия, Эдинбург.

[ПОКАЗАТЬ НА КАРТЕ](#)

Какая здесь погода?

Сейчас в Эдинбурге

+6 ° 

ясно от +6 до +9

Утро

+6



День

+8



Вечер

+4



А кто такая рапунцель?

Рапунцель



Сказка о девушке с очень длинными волосами, которая была заточена в высокой башне. Была записана братьями Гримм. В системе классификации народных сказок Аарне-Томпсона этот сюжет имеет номер 310: «Дева в башне».

ru.wikipedia.org

[ОТКРЫТЬ В ПОИСКЕ](#)

Да это же кружка с рапунцией.

[НАЙТИ ЭТО В ЯНДЕКСЕ](#)

[ПОСМОТРЕТЬ ПОХОЖИЕ КАРТИНКИ](#)

An example of a dialogue system

Что ты делаешь?

Новости тут почитала. Если коротко, то мир меняется, стабильности нет.

А я лекцию готовлю.

Ха, и как успехи?

Да вот пытаюсь найти примеры диалоговых систем.

Что ж, удачи вам в этом нелегком деле

Алиса, ну, помоги мне подготовить лекцию.

Чем смогу, помогу

Без тебя никак.

Вы даже не представляете, как приятно это слышать

Лекция про тебя будет.

О, я надеюсь на это

Task-oriented

Goal: solve a task for a user



General conversation

Goal: conversation itself

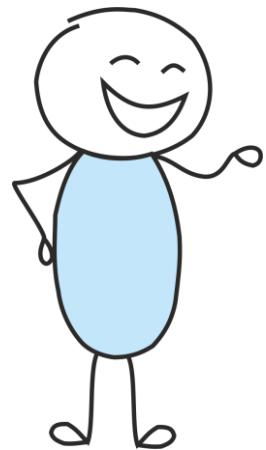
It has to be enjoyable!



Task-oriented

AKA GOAL-ORIENTED

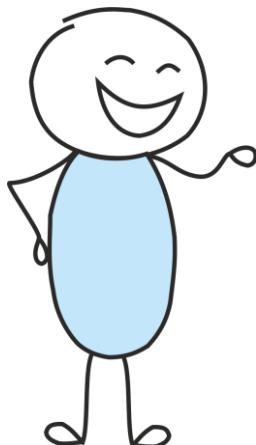
Task-oriented dialogue system



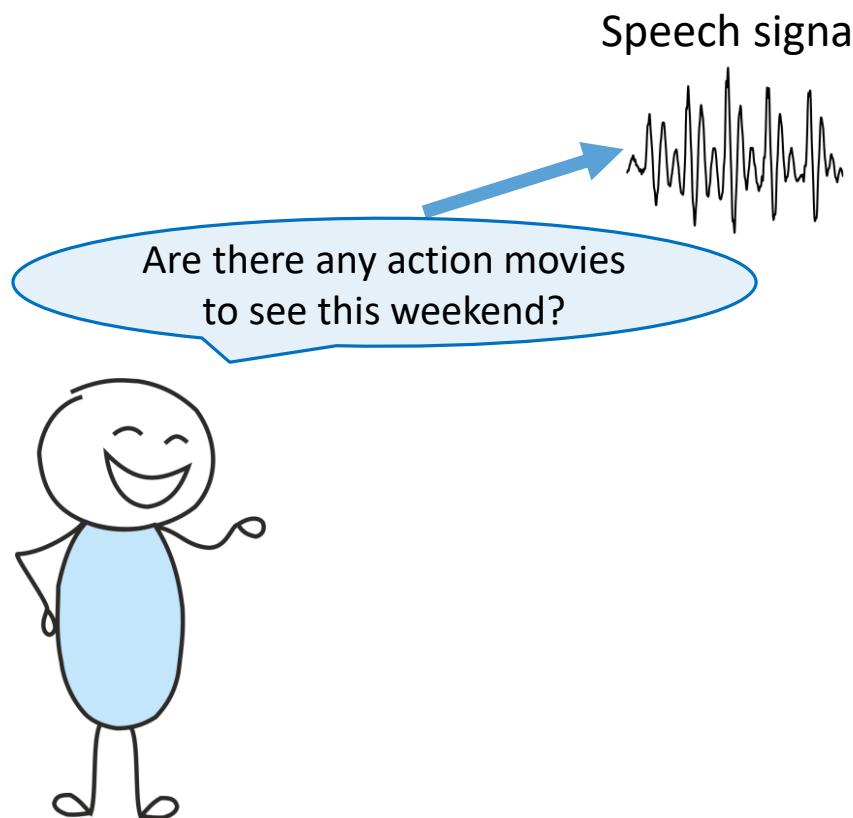
Task-oriented dialogue system



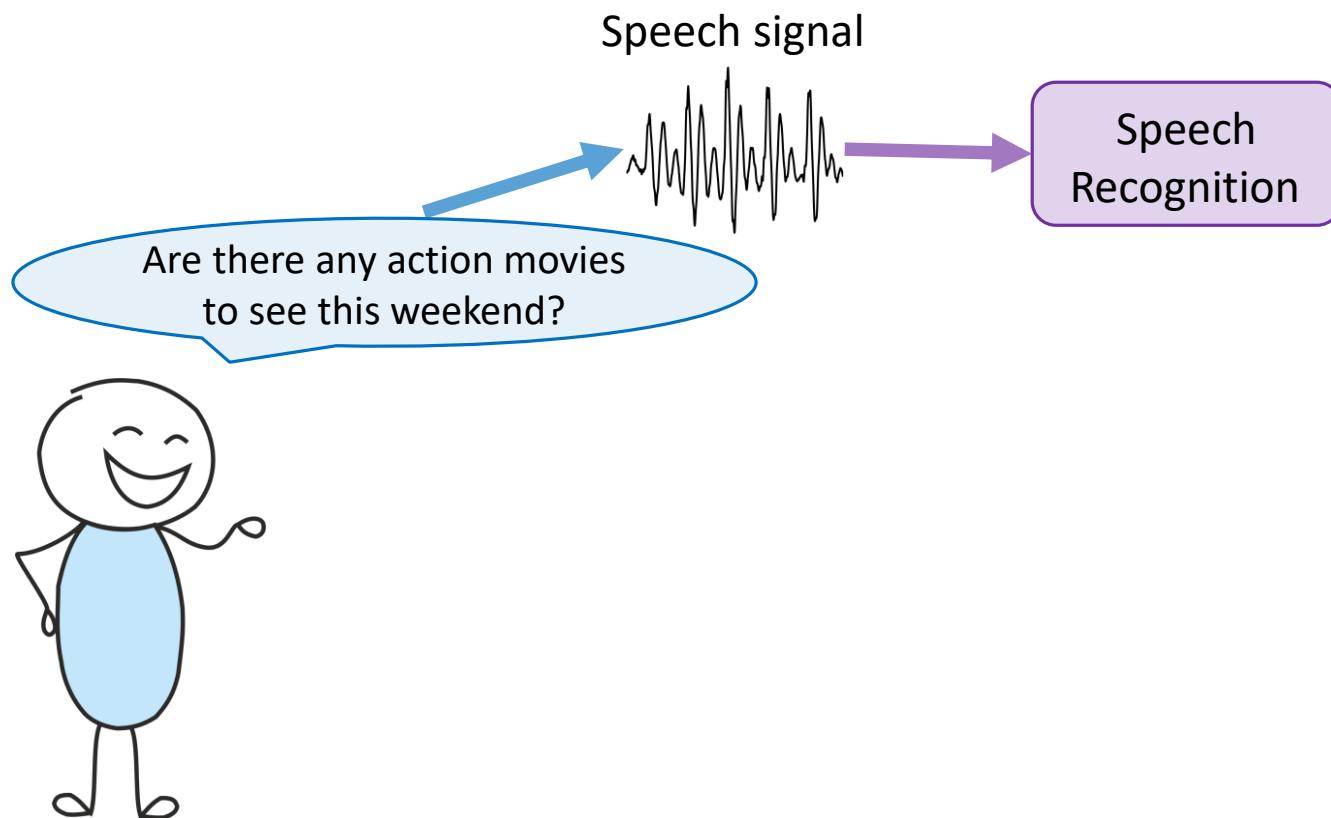
Are there any action movies
to see this weekend?



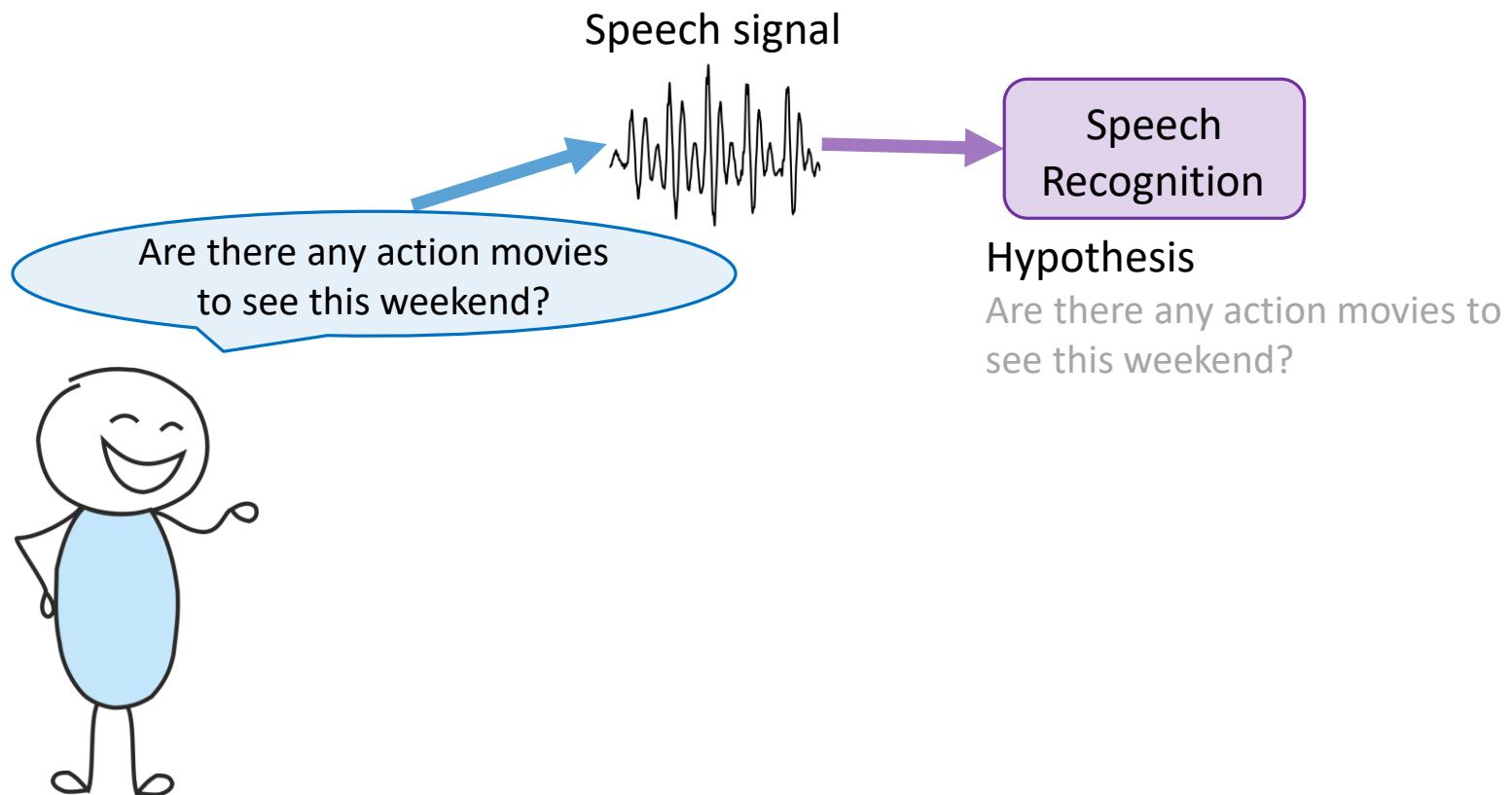
Task-oriented dialogue system



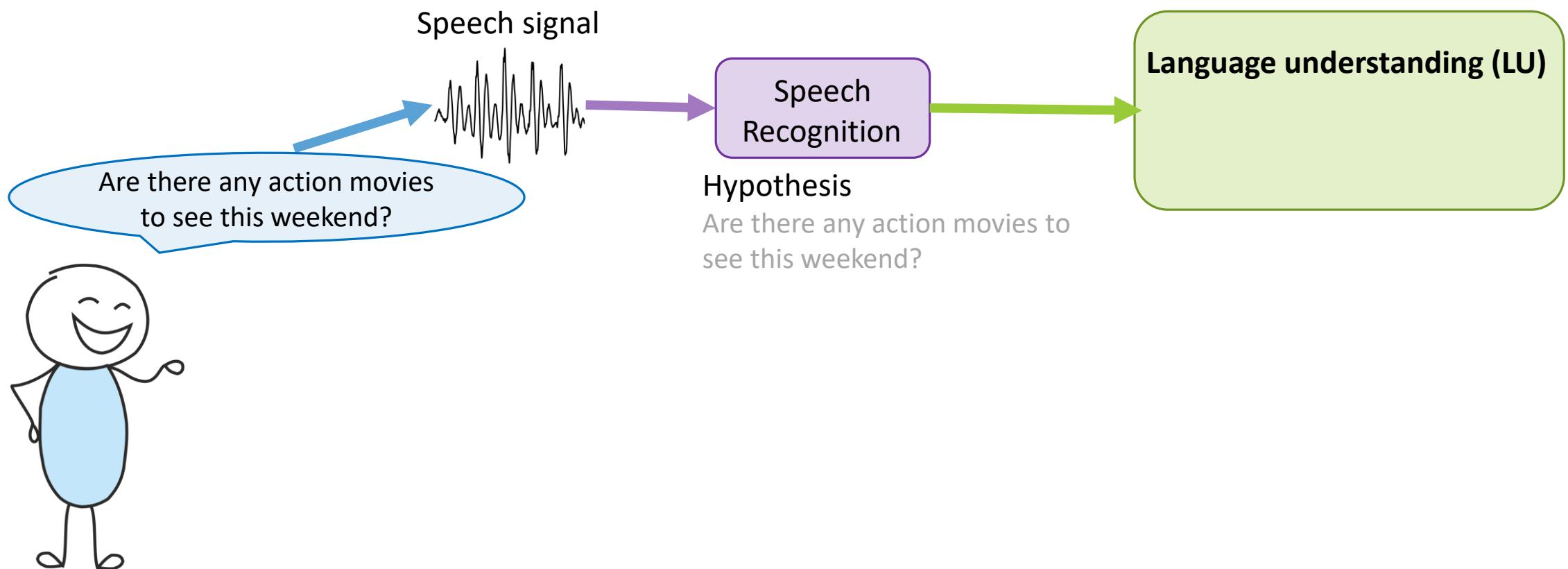
Task-oriented dialogue system



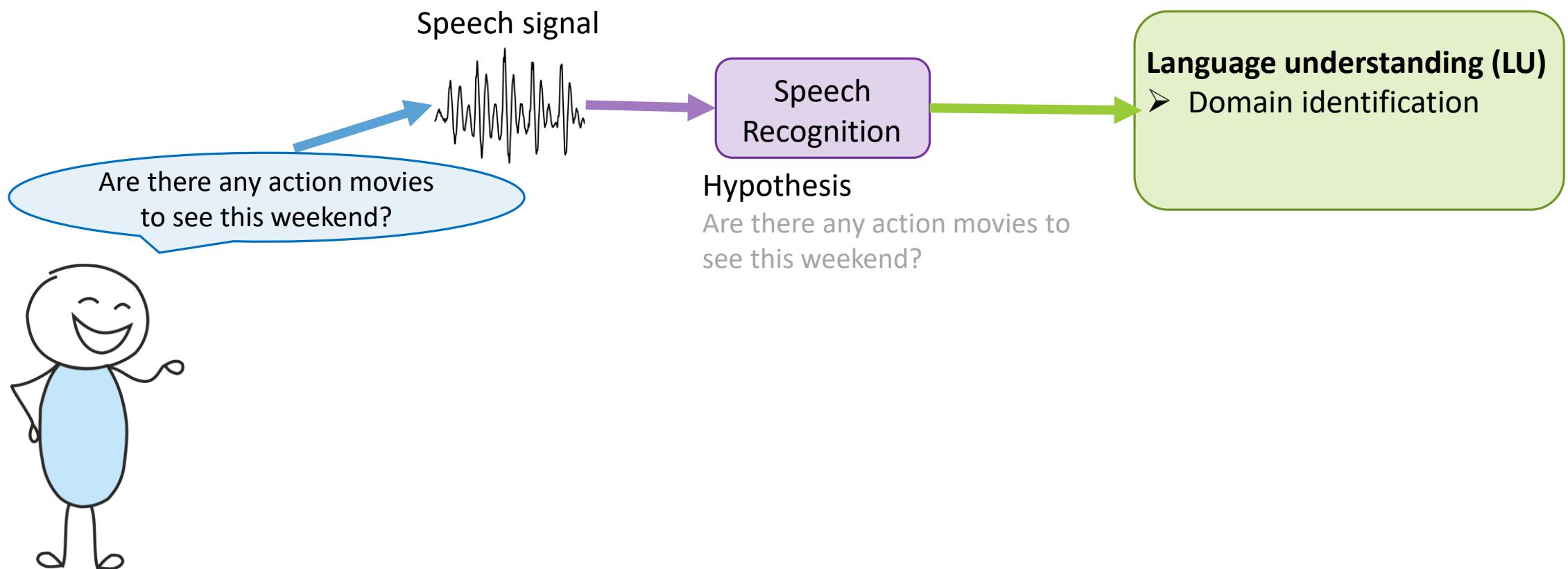
Task-oriented dialogue system



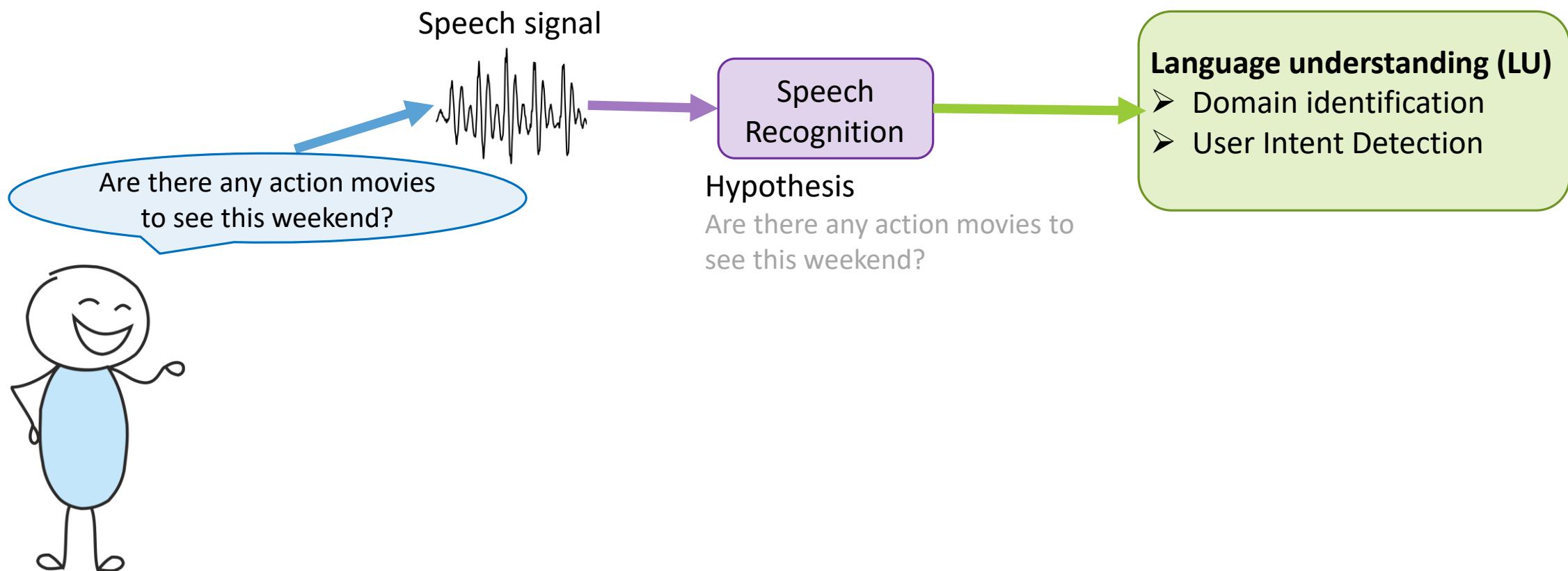
Task-oriented dialogue system



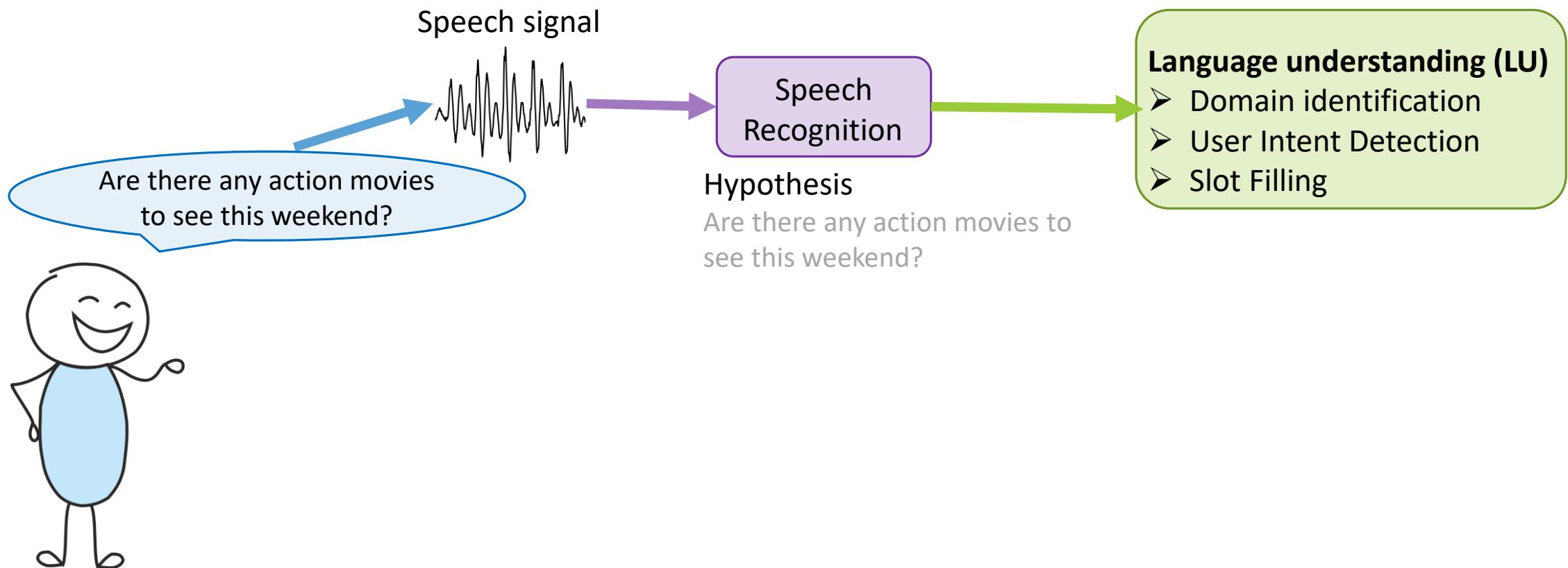
Task-oriented dialogue system



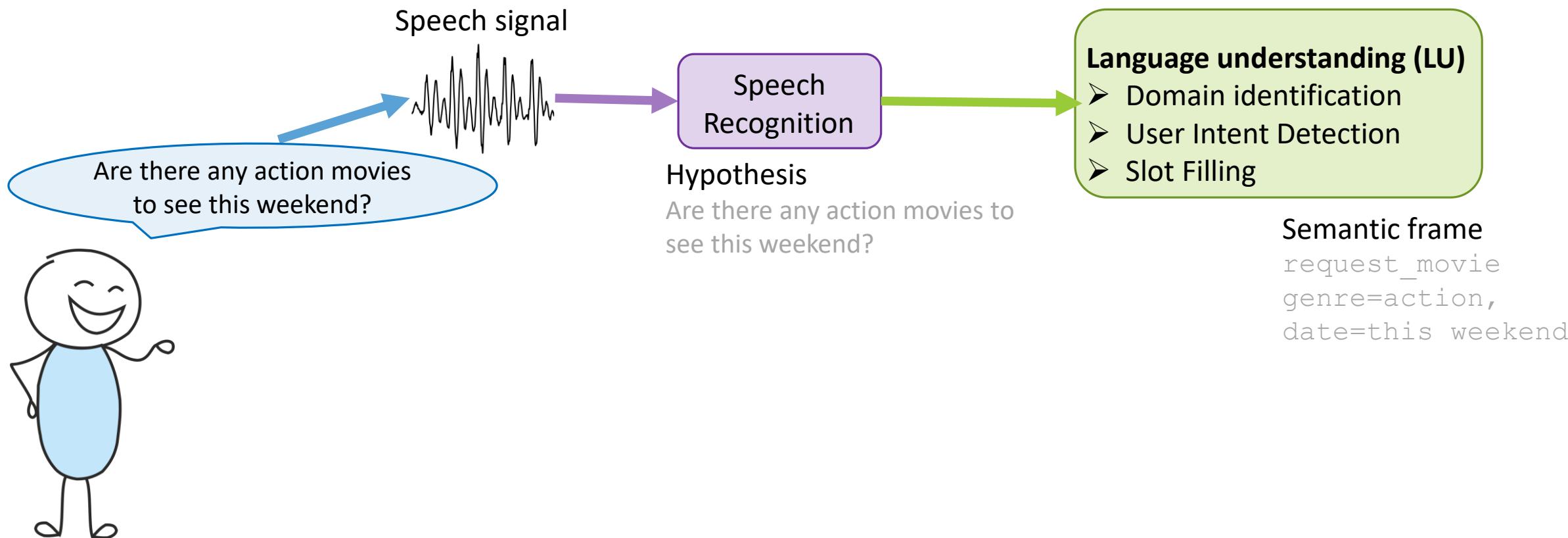
Task-oriented dialogue system



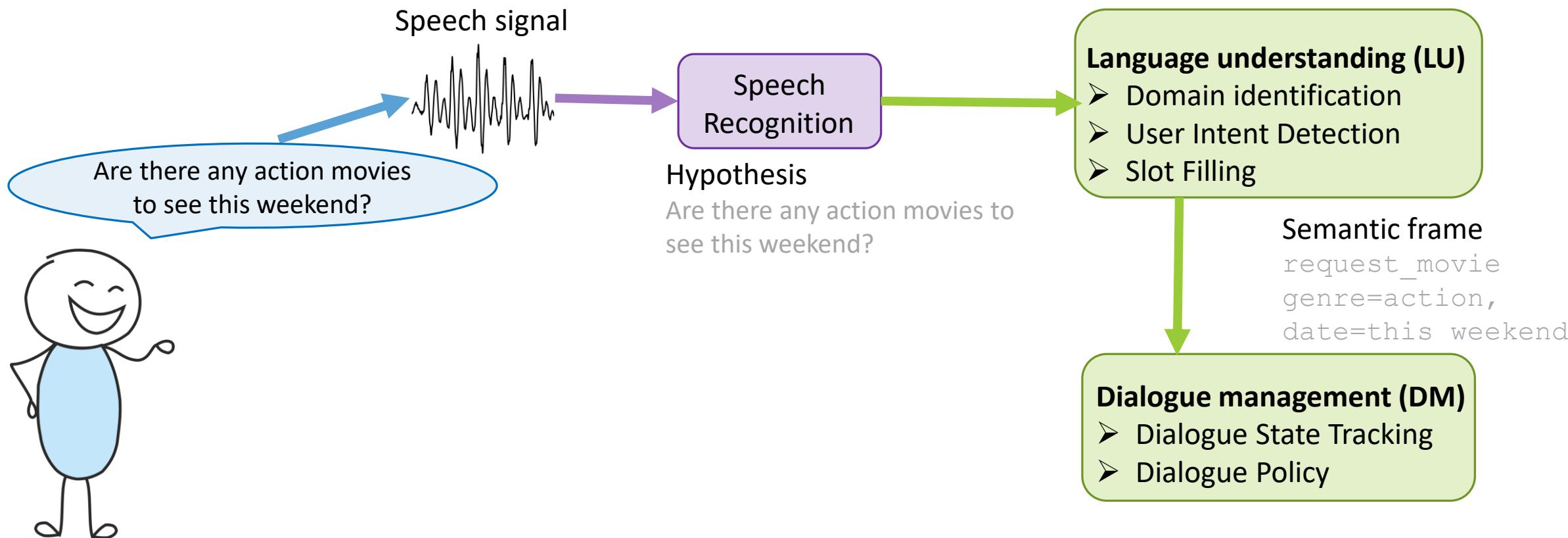
Task-oriented dialogue system



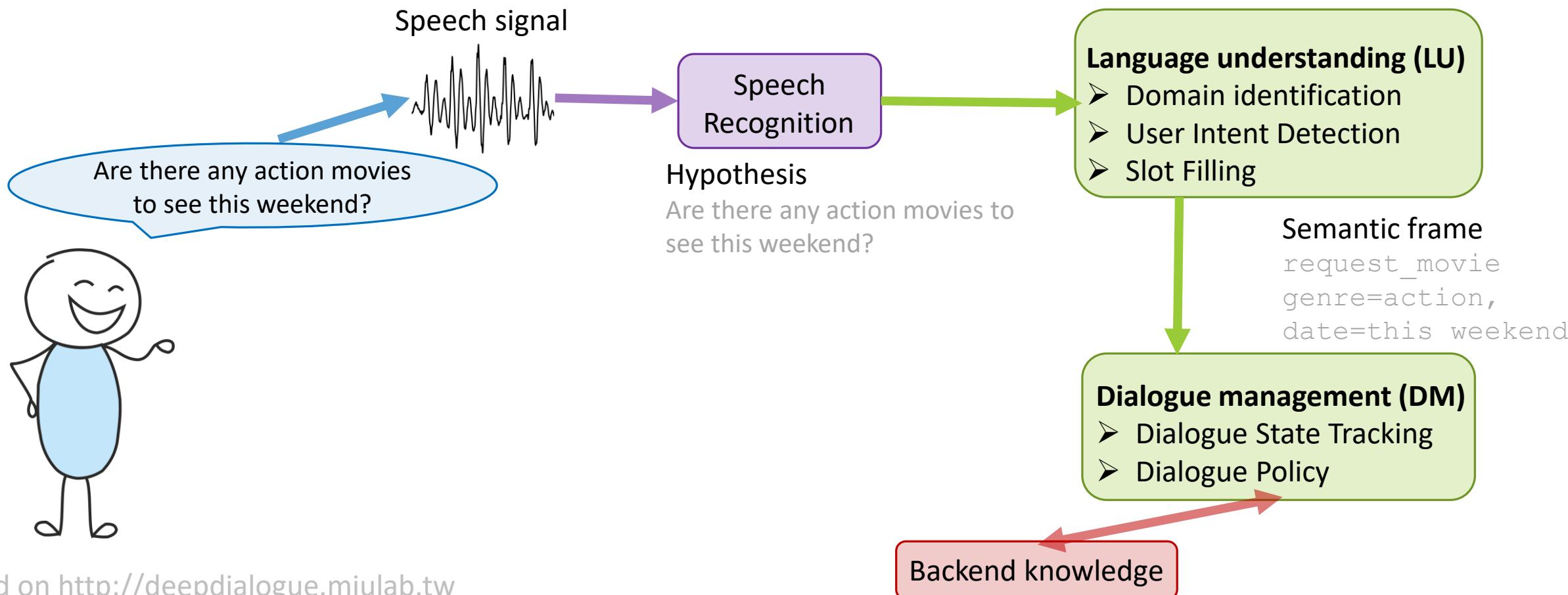
Task-oriented dialogue system



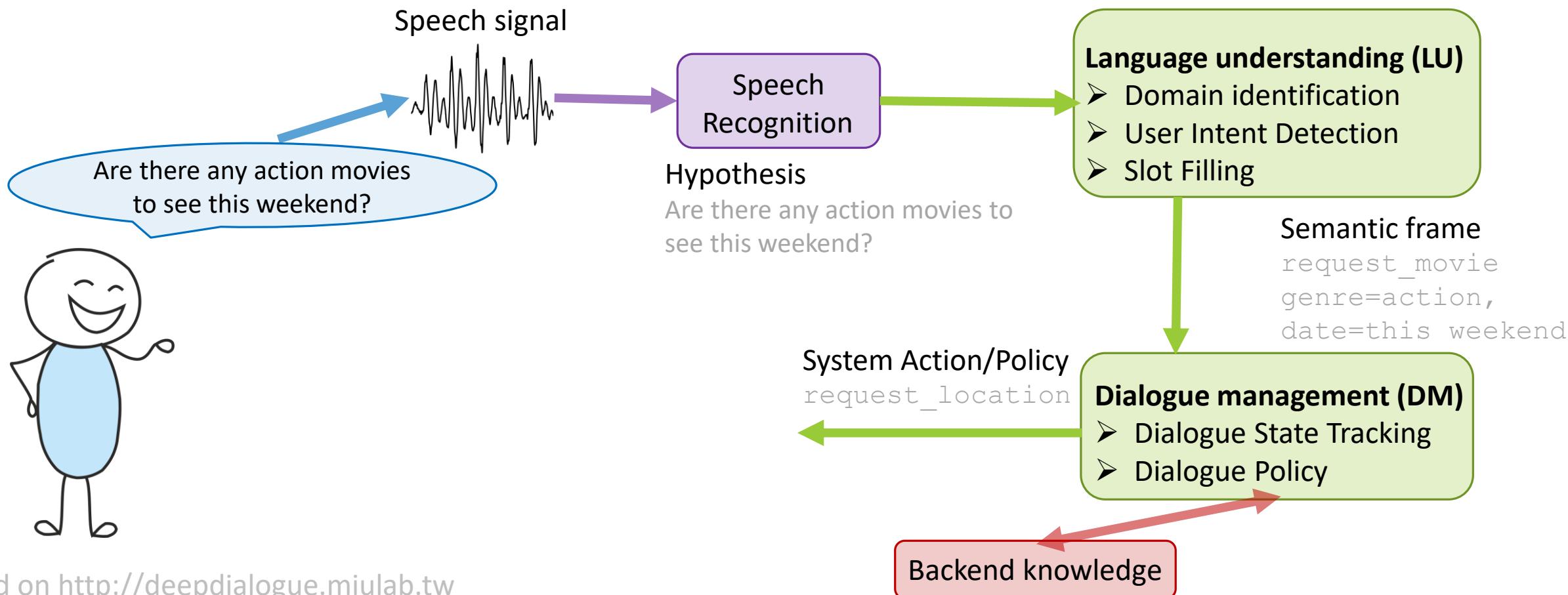
Task-oriented dialogue system



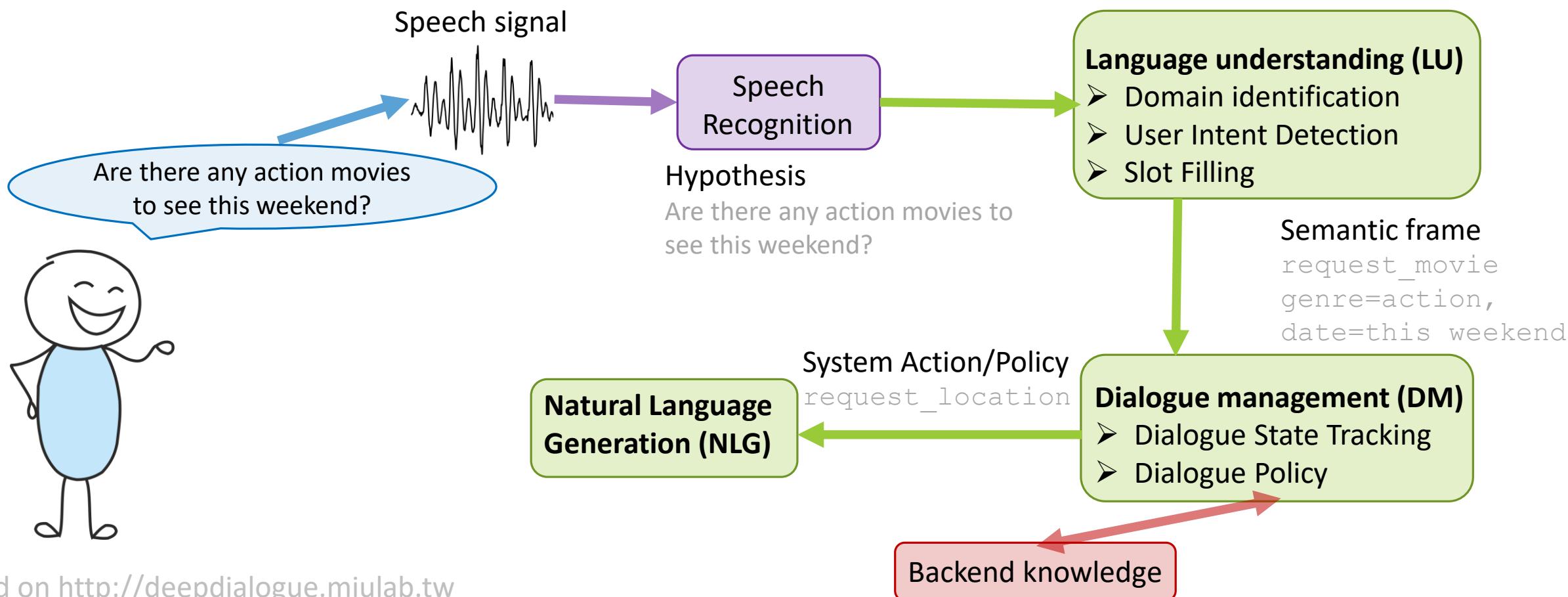
Task-oriented dialogue system



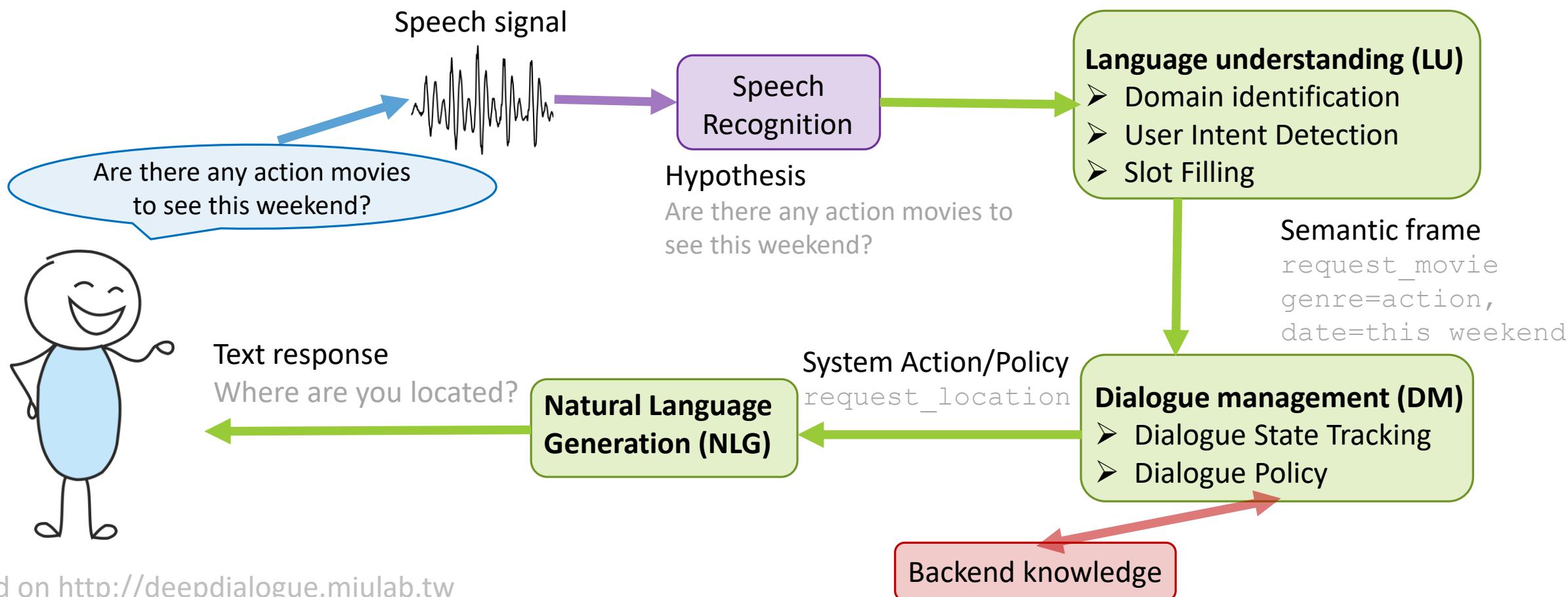
Task-oriented dialogue system



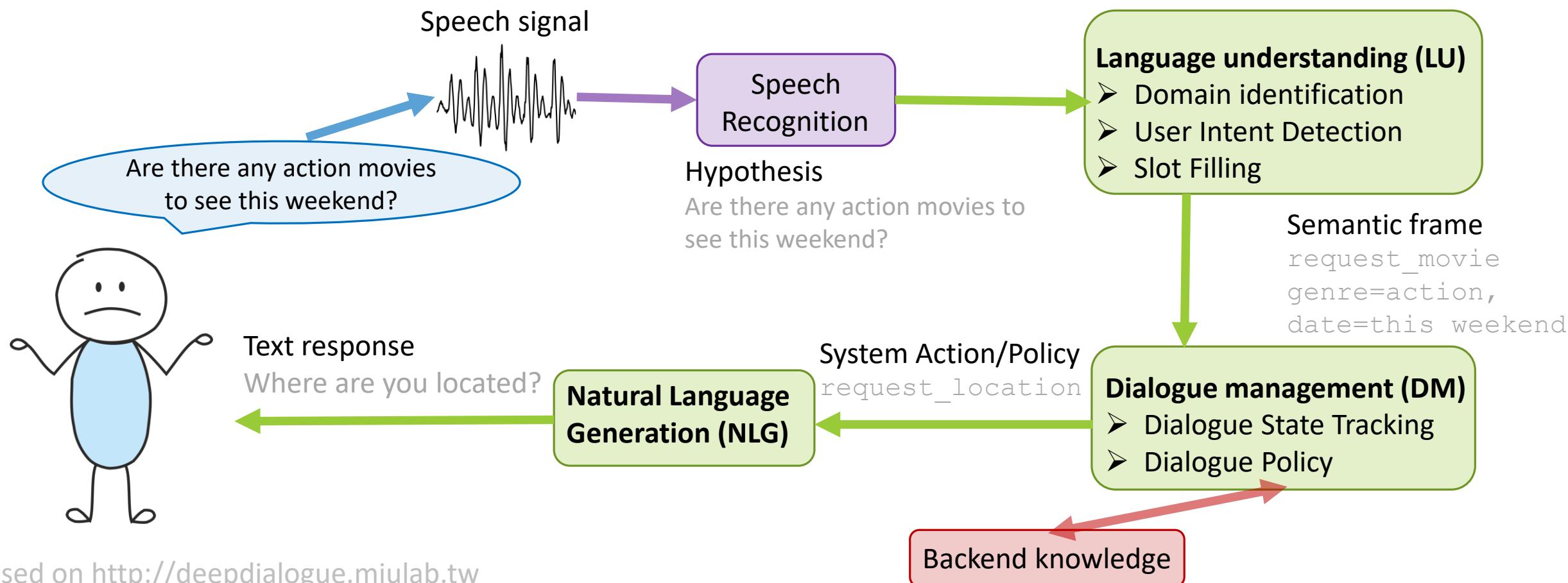
Task-oriented dialogue system



Task-oriented dialogue system

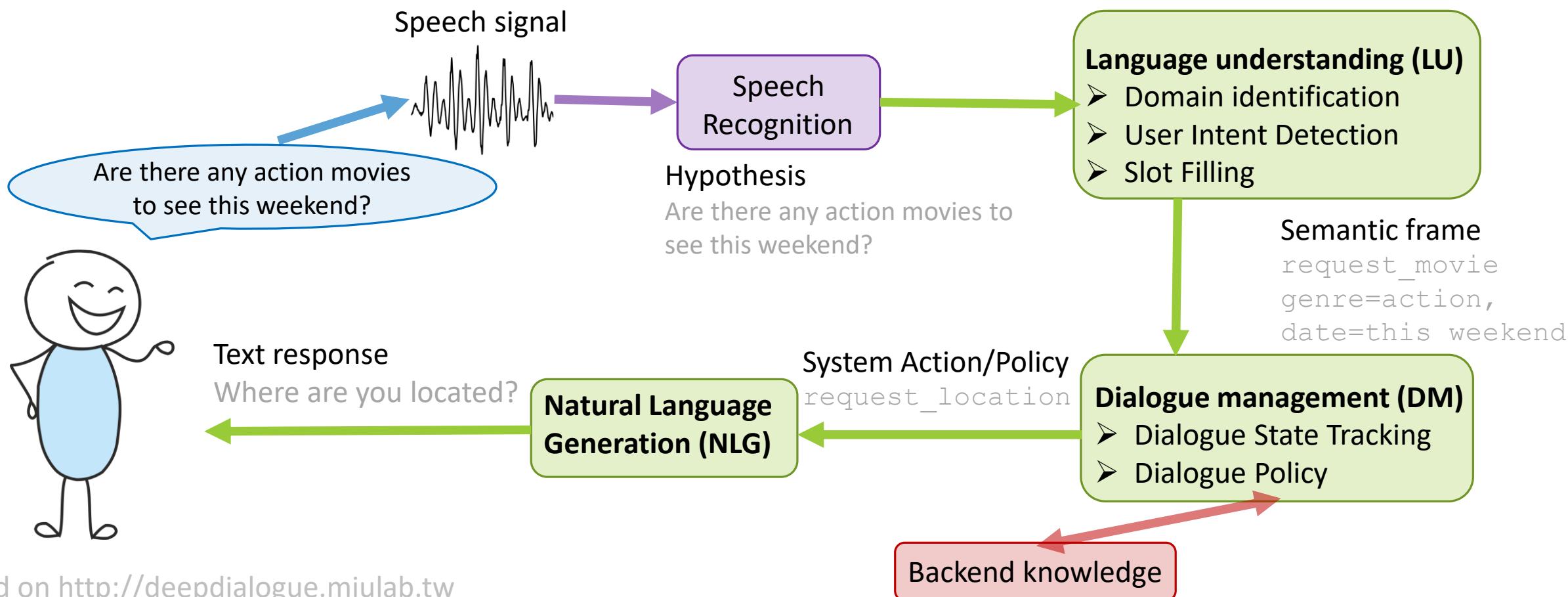


Task-oriented dialogue system

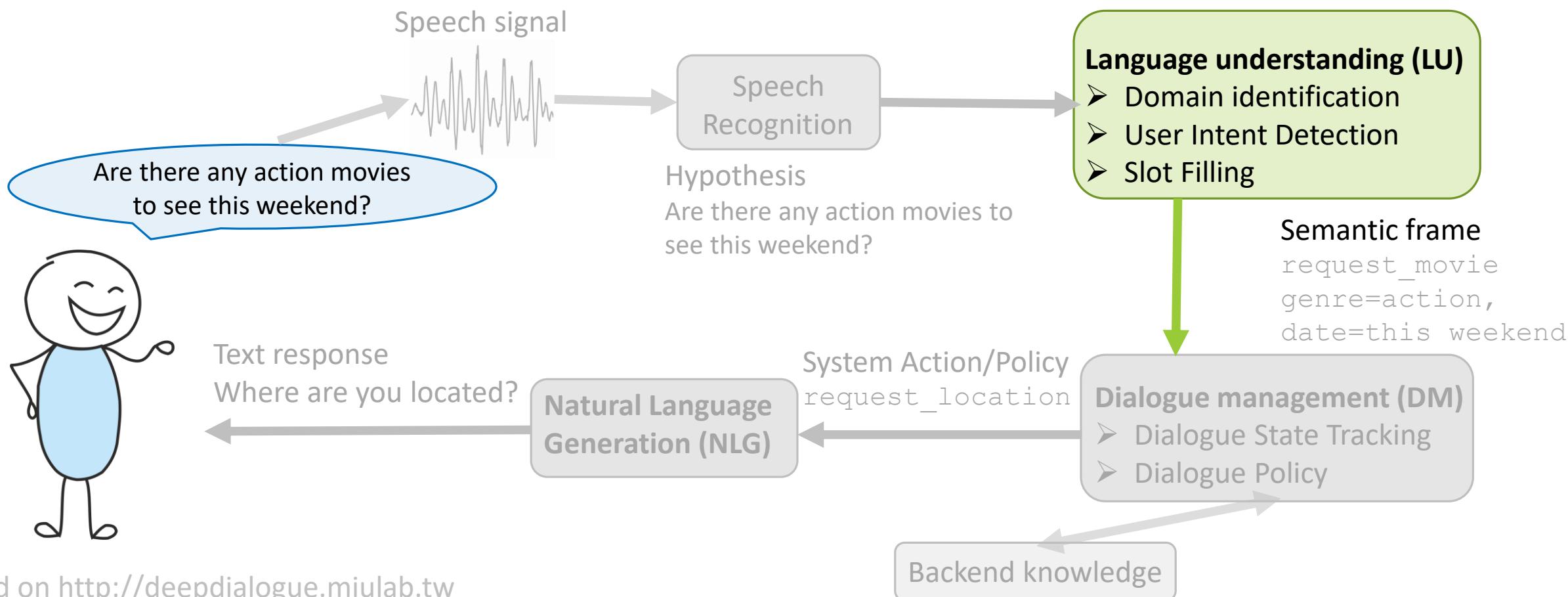


Language Understanding

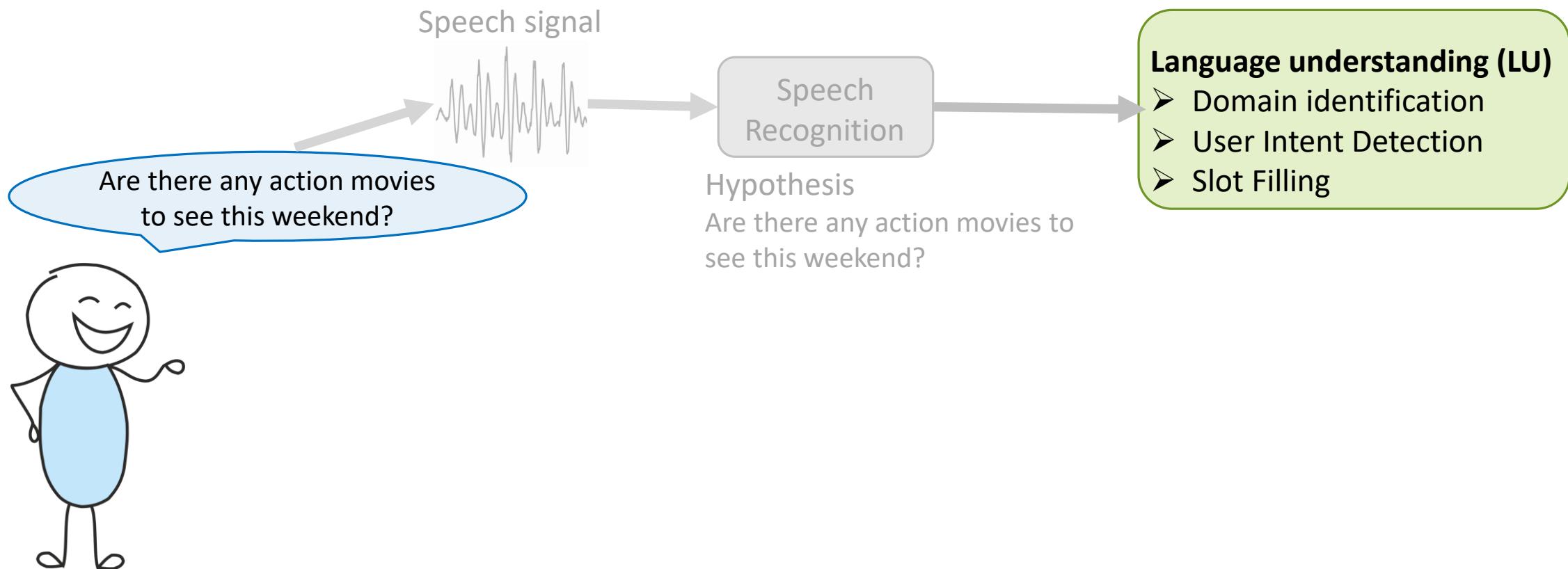
Task-oriented dialogue system



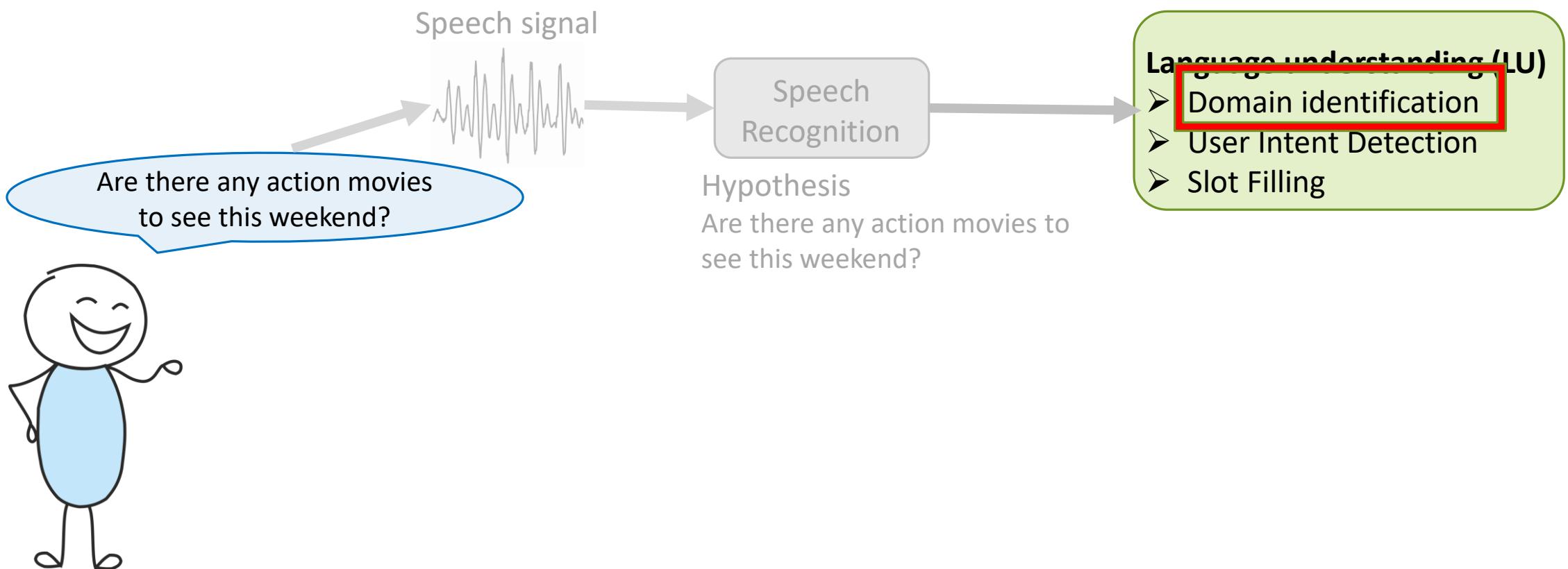
Language Understanding



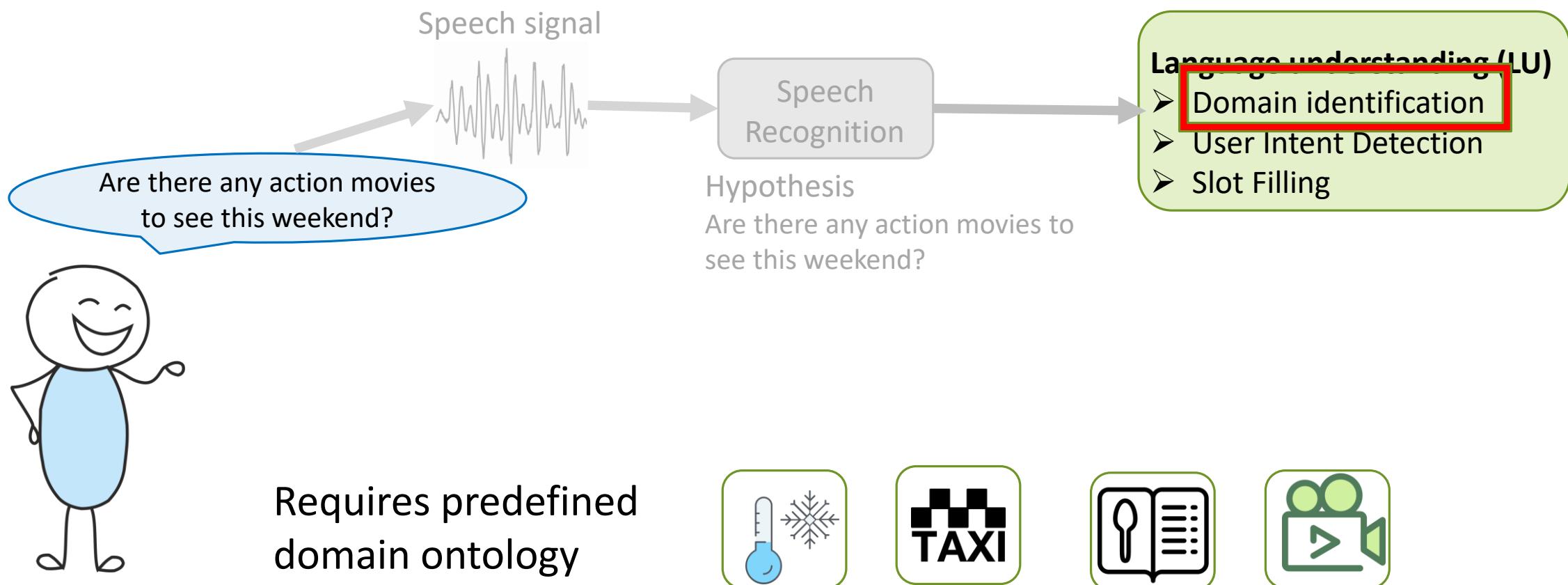
Language Understanding



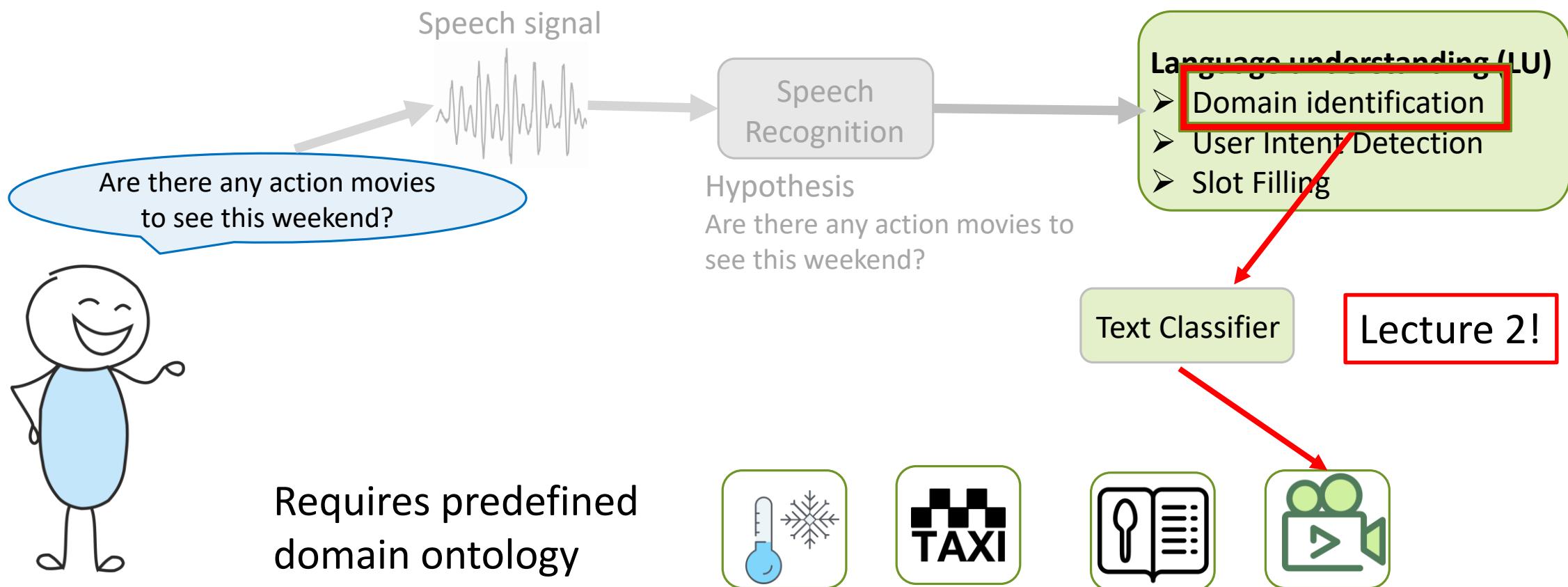
Domain Identification



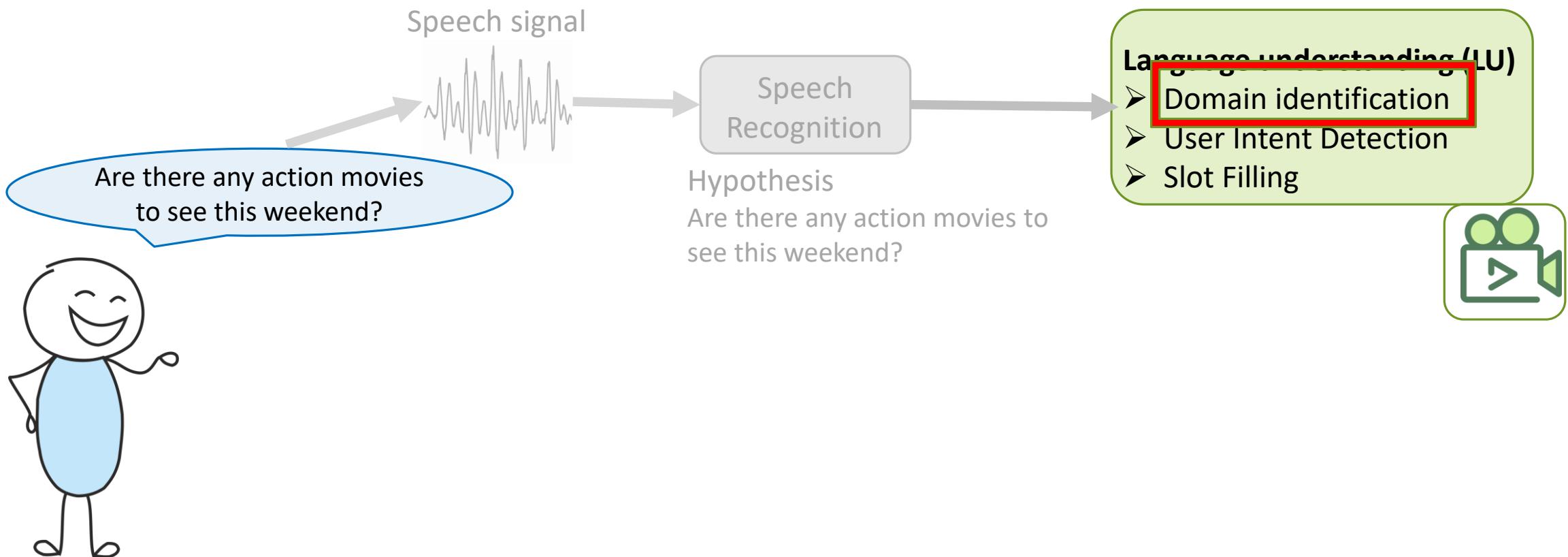
Domain Identification



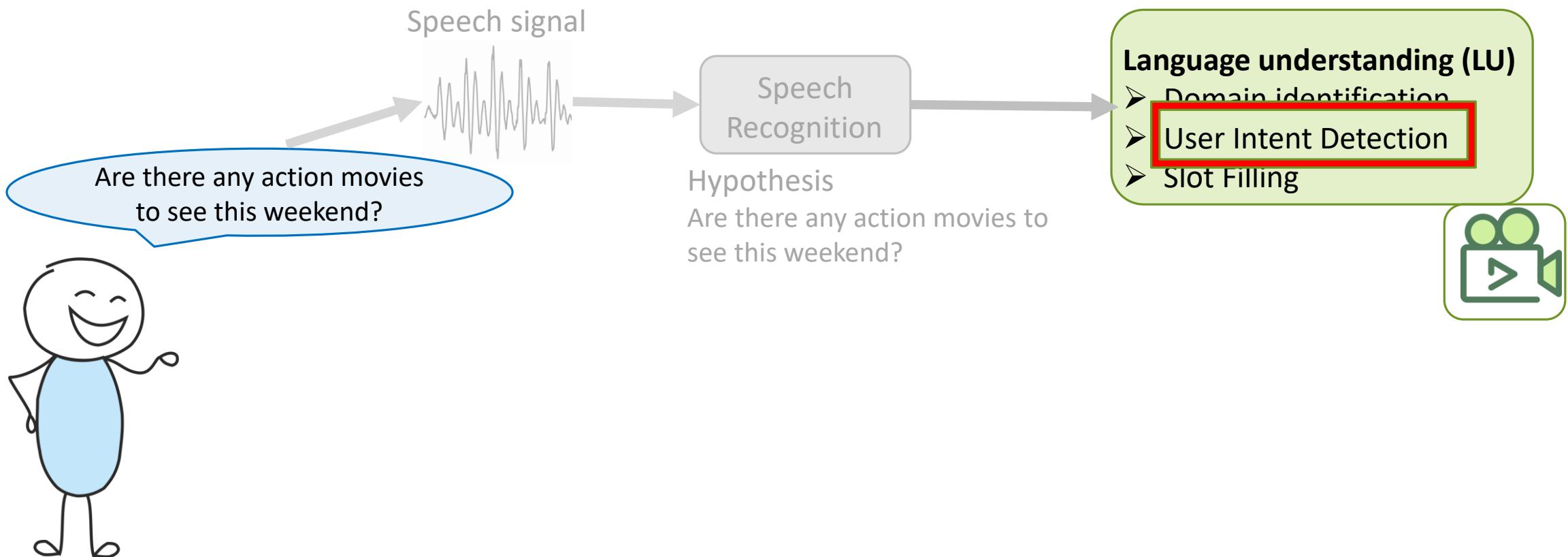
Domain Identification



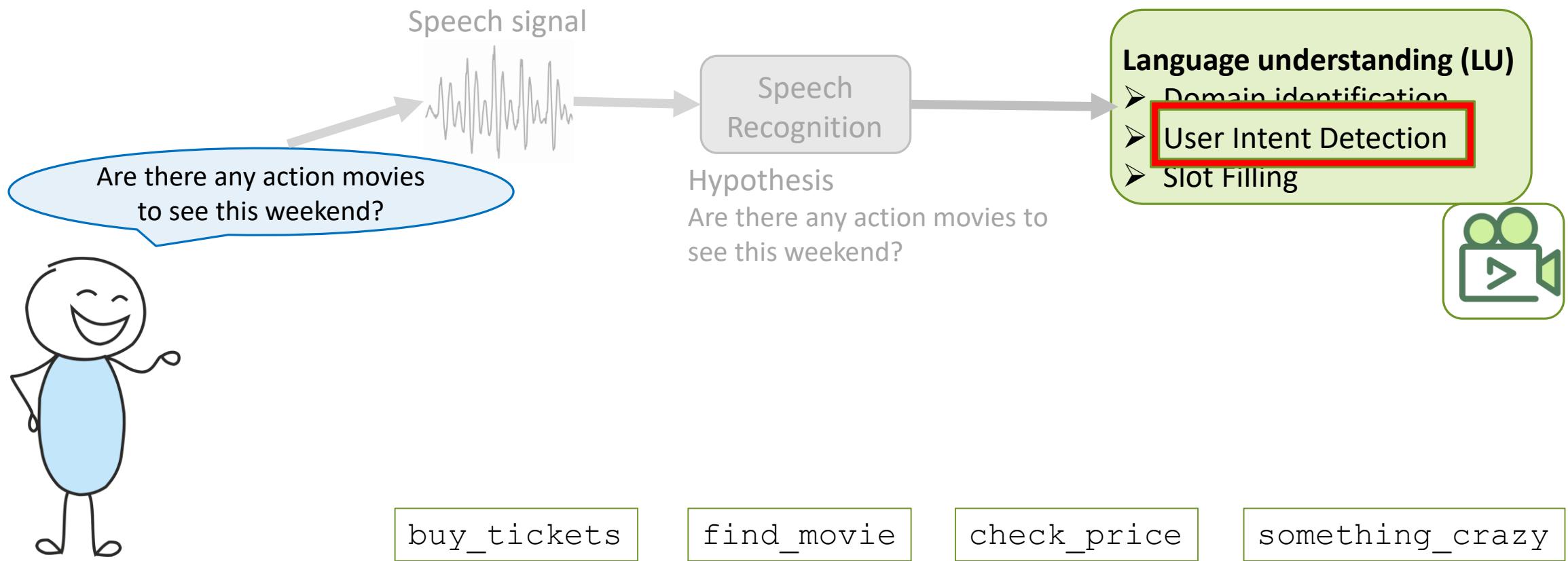
Domain Identification



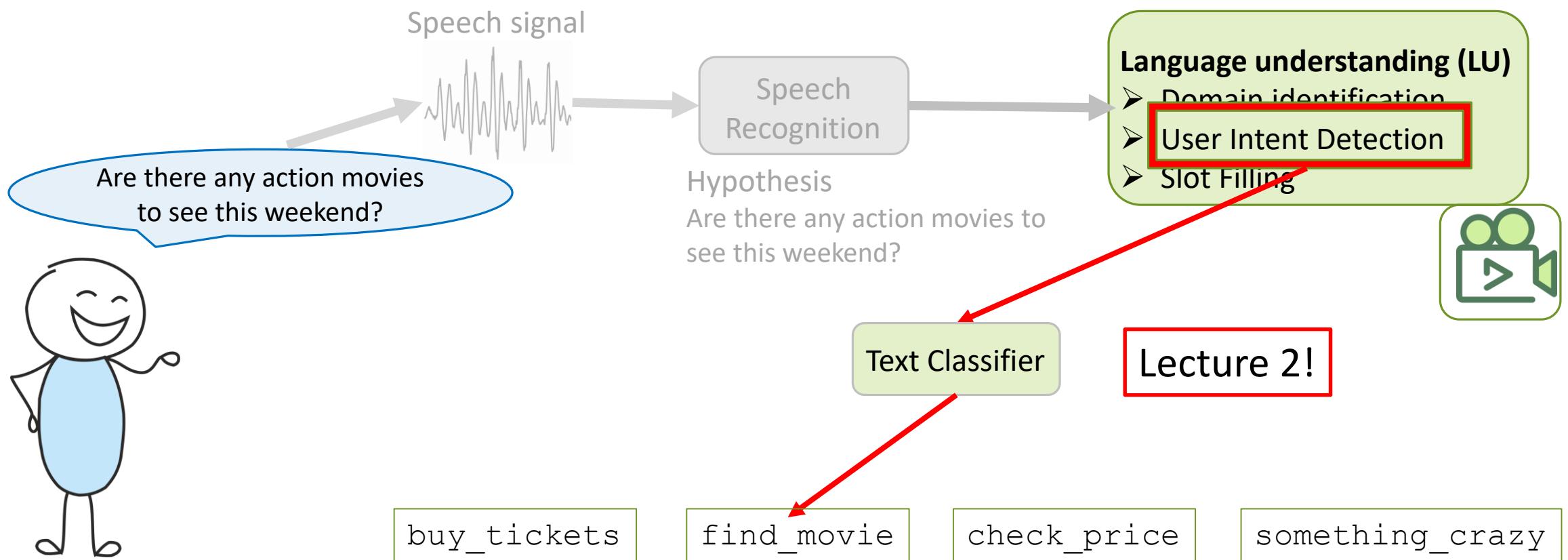
User Intent Detection



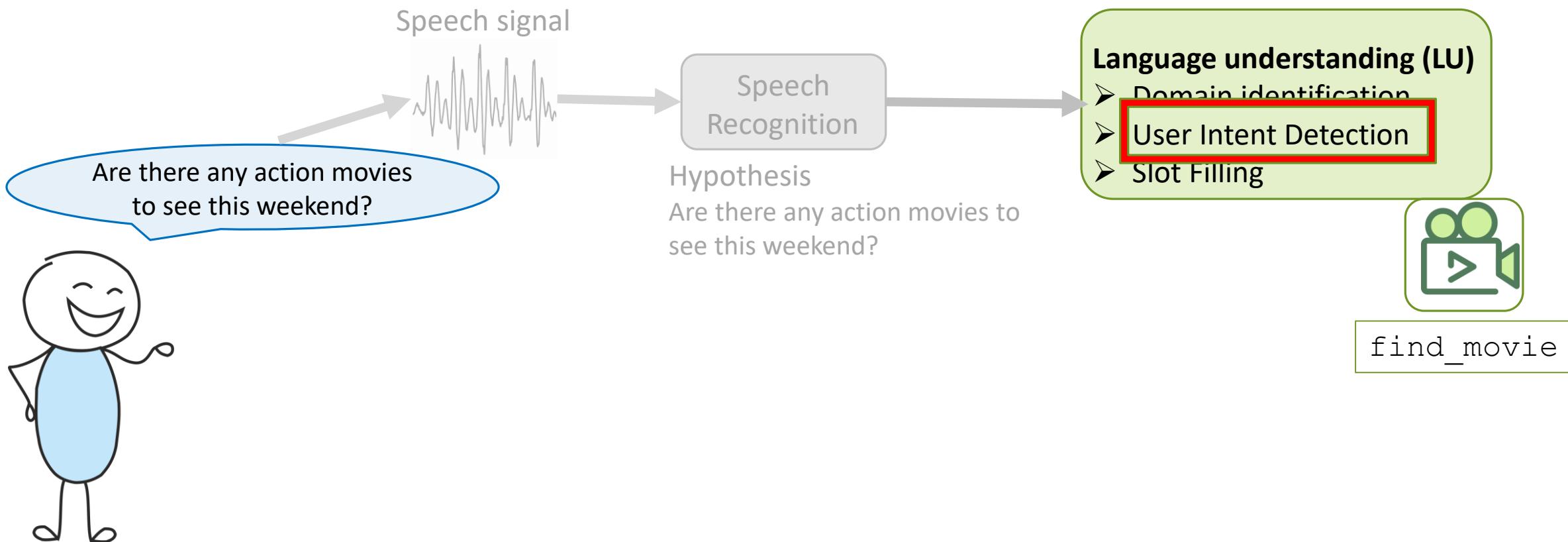
User Intent Detection



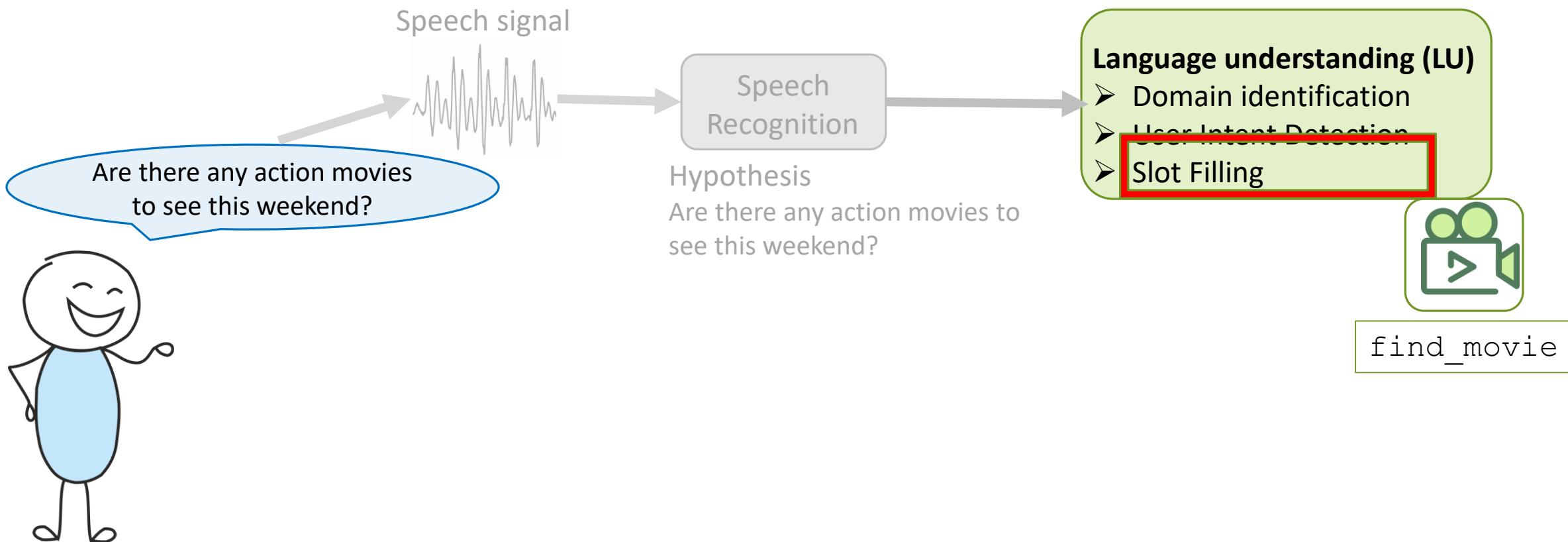
User Intent Detection



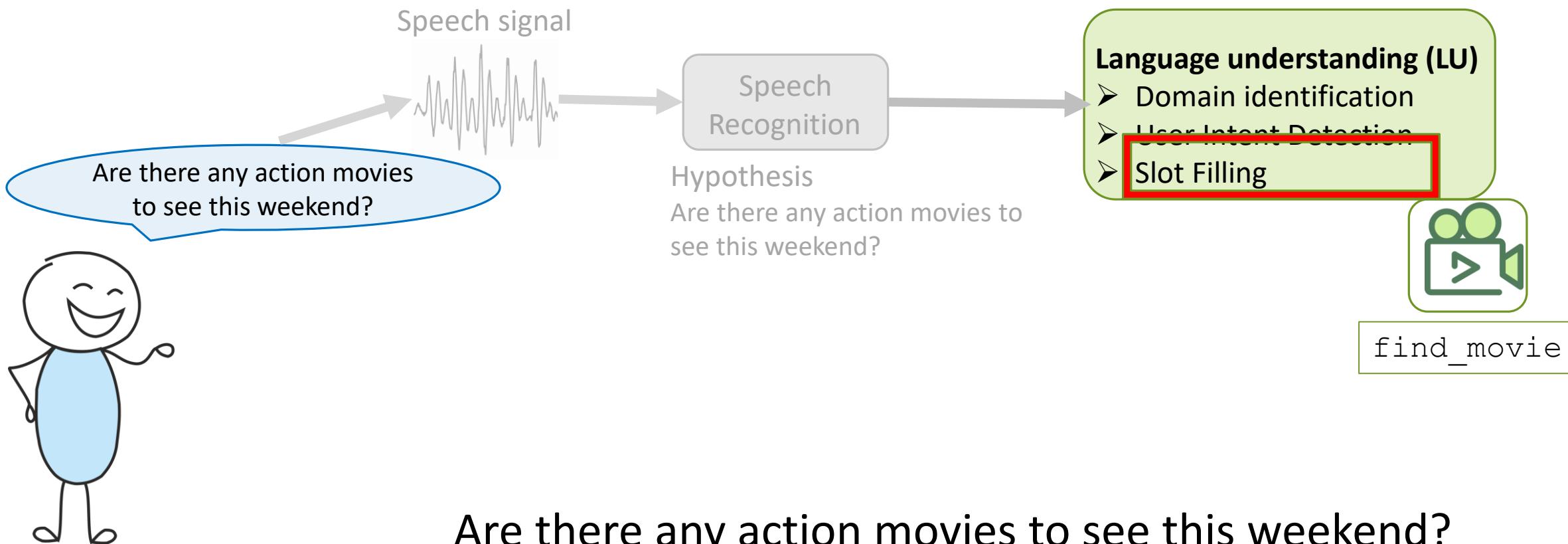
User Intent Detection



Slot Filling

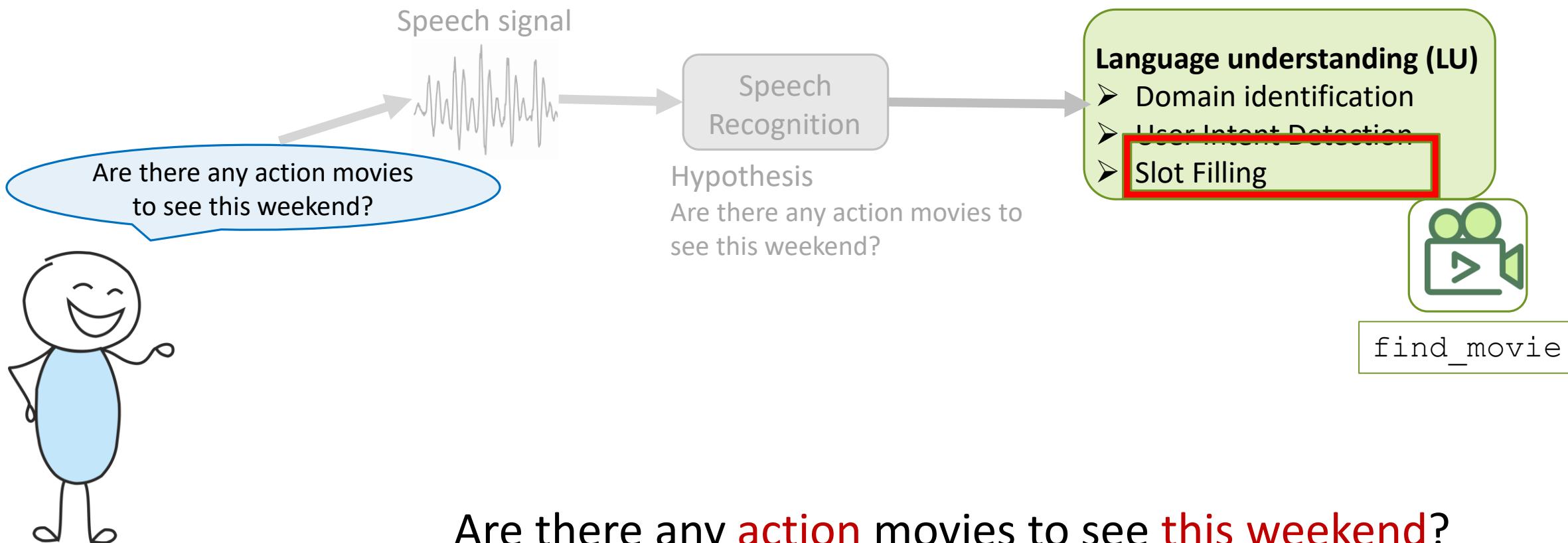


Slot Filling



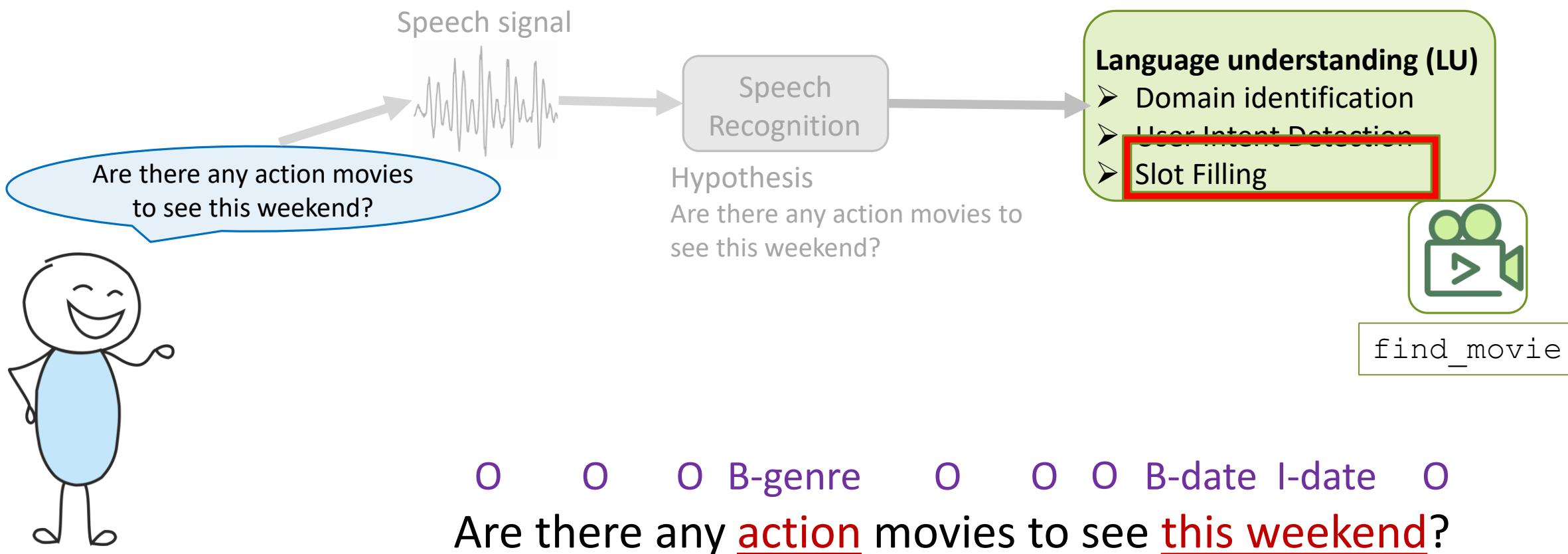
Are there any action movies to see this weekend?

Slot Filling

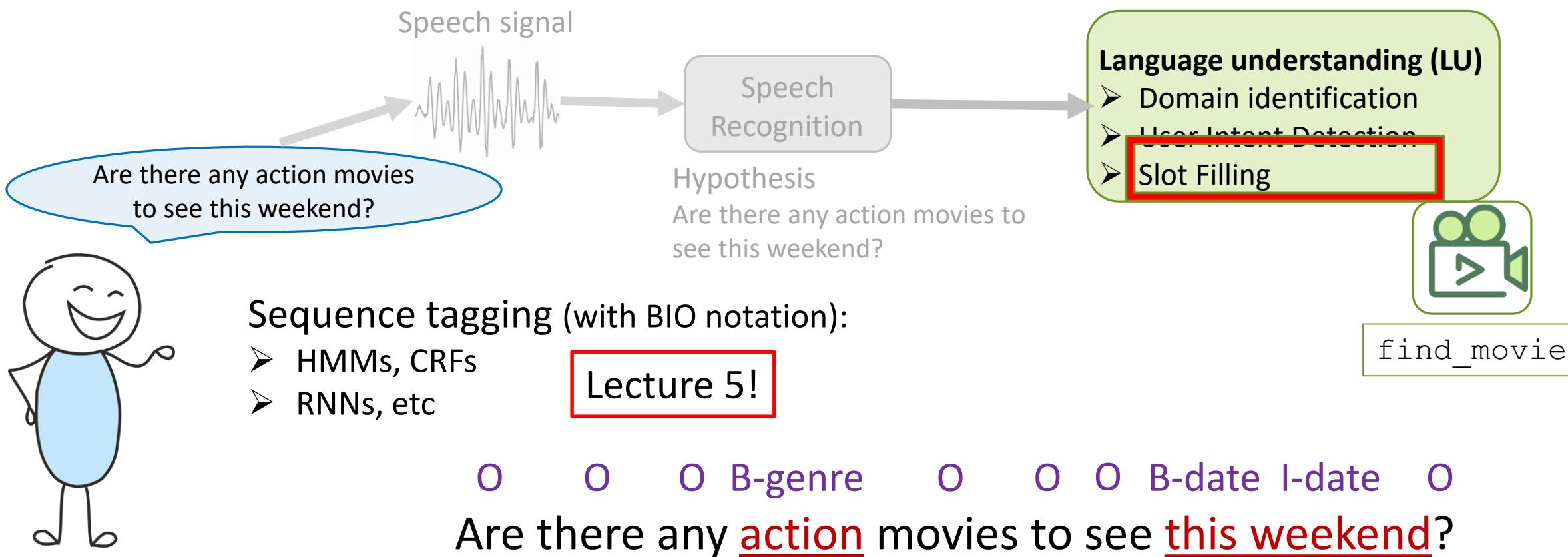


Are there any action movies to see this weekend?

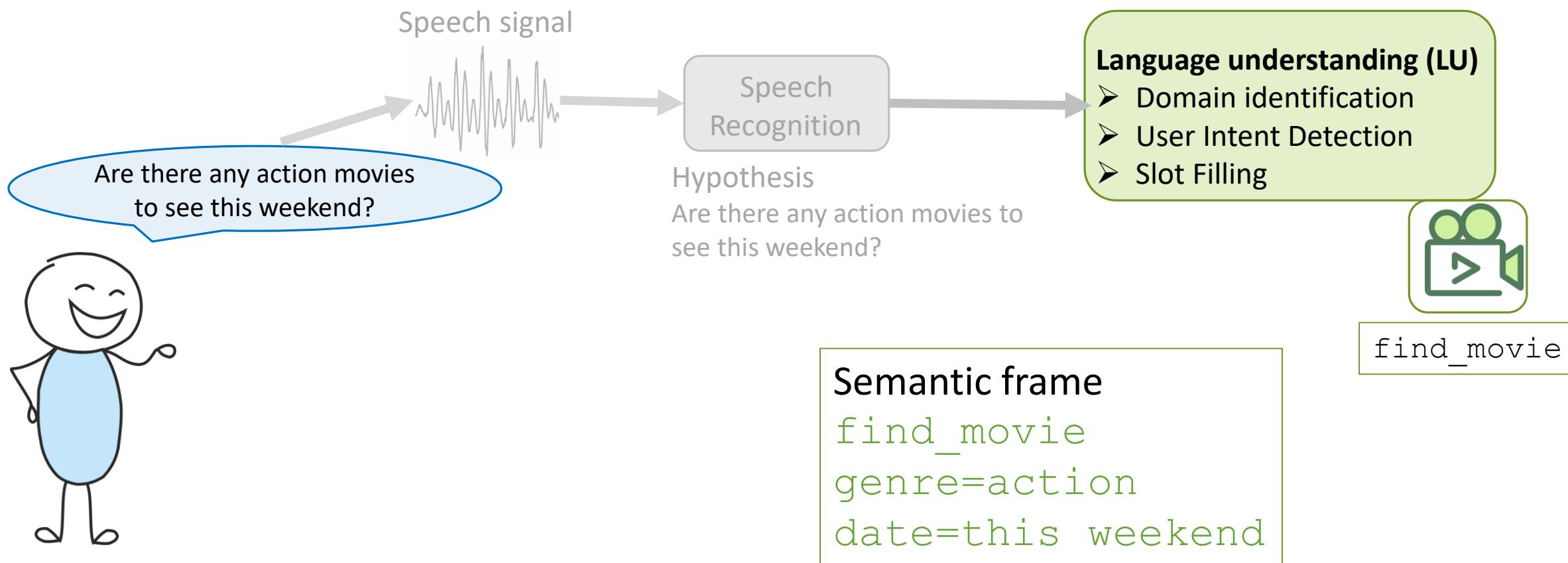
Slot Filling



Slot Filling

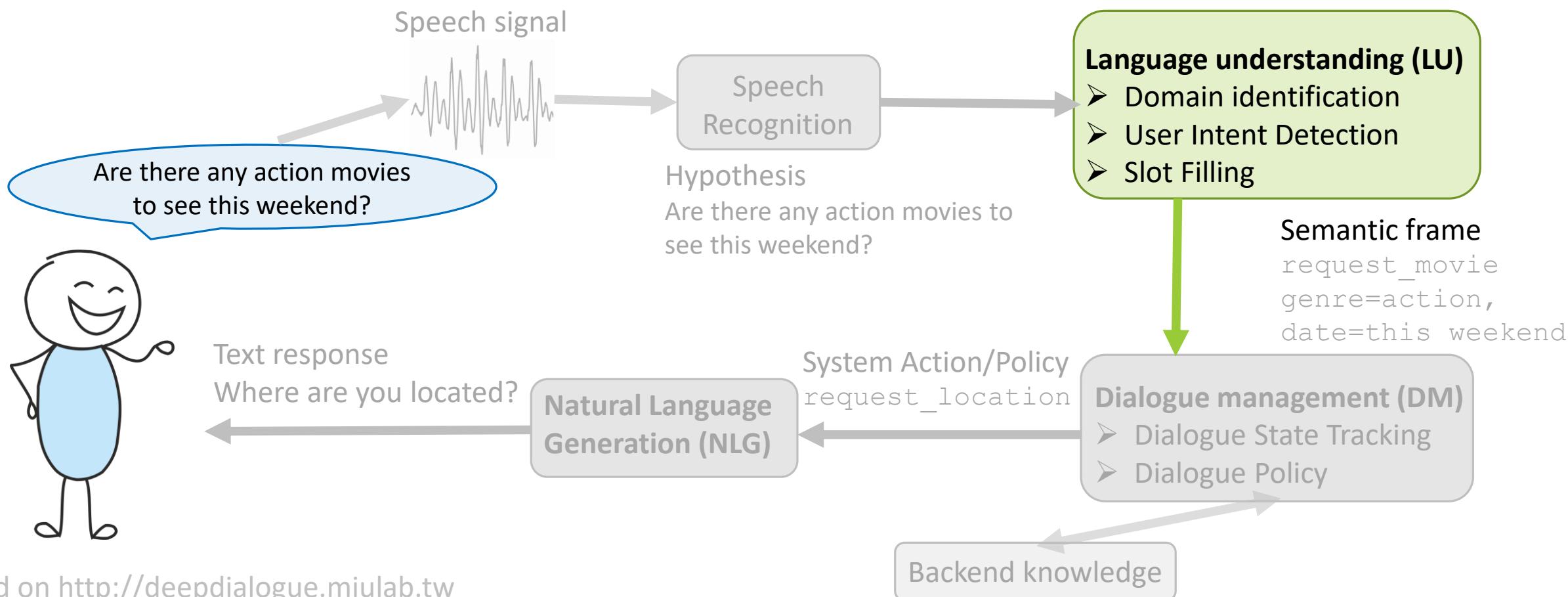


Language Understanding

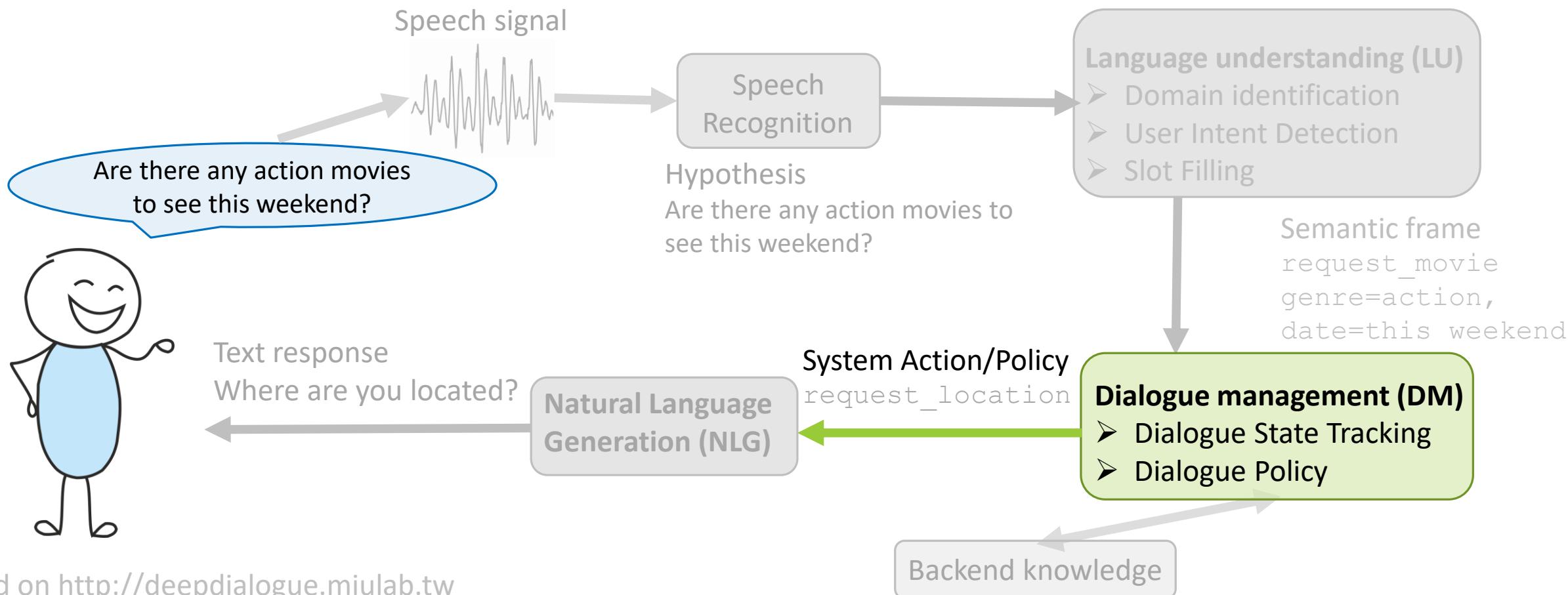


Dialogue Management

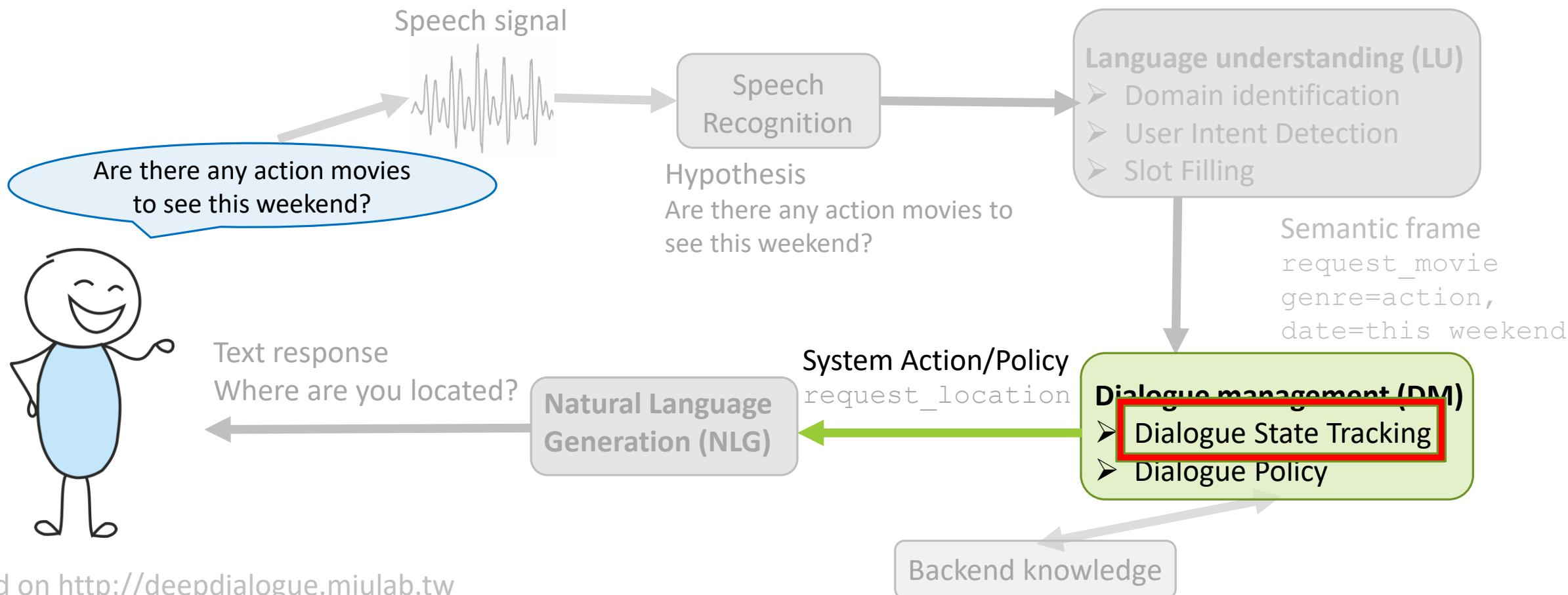
Language Understanding



Dialogue Management



Dialogue Management



Dialogue State

- goal constraint – constraints from a user
 - Some value
 - “Dontcare” – no preference
 - None – need to be specified
- requested slots – slots that are to be specified
- search method – how a user interacts with a system
 - constraints
 - alternatives
 - finished

| | SLU output | Dialog State |
|--------|----------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------|
| Turn 1 | | |
| a: | What part of town did you have in mind? <i>request(area)</i> | 0.7 <i>null()</i> 0.2 <i>inform(food='north african')</i> 0.1 <i>inform(area=north)</i> |
| u: | The North area <i>inform(area=north)</i> | goal constraints <i>area=north</i> requested slots – search method <i>by constraints</i> |
| Turn 2 | | |
| a: | Which part of town? <i>request(area)</i> | 0.8 <i>inform(area=north, pricerange=cheap)</i> 0.2 <i>inform(area=north)</i> |
| u: | A cheap place in the North <i>inform(area=north, pricerange=cheap)</i> | goal constraints <i>area=north</i> requested slots – search method <i>by constraints</i> |

Dialogue State

- goal constraint – constraints from a user
 - Some value
 - “Dontcare” – no preference
 - None – need to be specified
- requested slots – slots that are to be specified
- search method – how a user interacts with a system
 - constraints
 - alternatives
 - finished

Turn 3

a: Da Vinci Pizzeria
is a cheap place in
the North.
*inform(name='Da
Vinci Pizzeria',
area=north,
pricerange=cheap)*

u: Do you have any-
thing else, but in
the West?
*requestAlternatives(
area=west)*

Turn 4

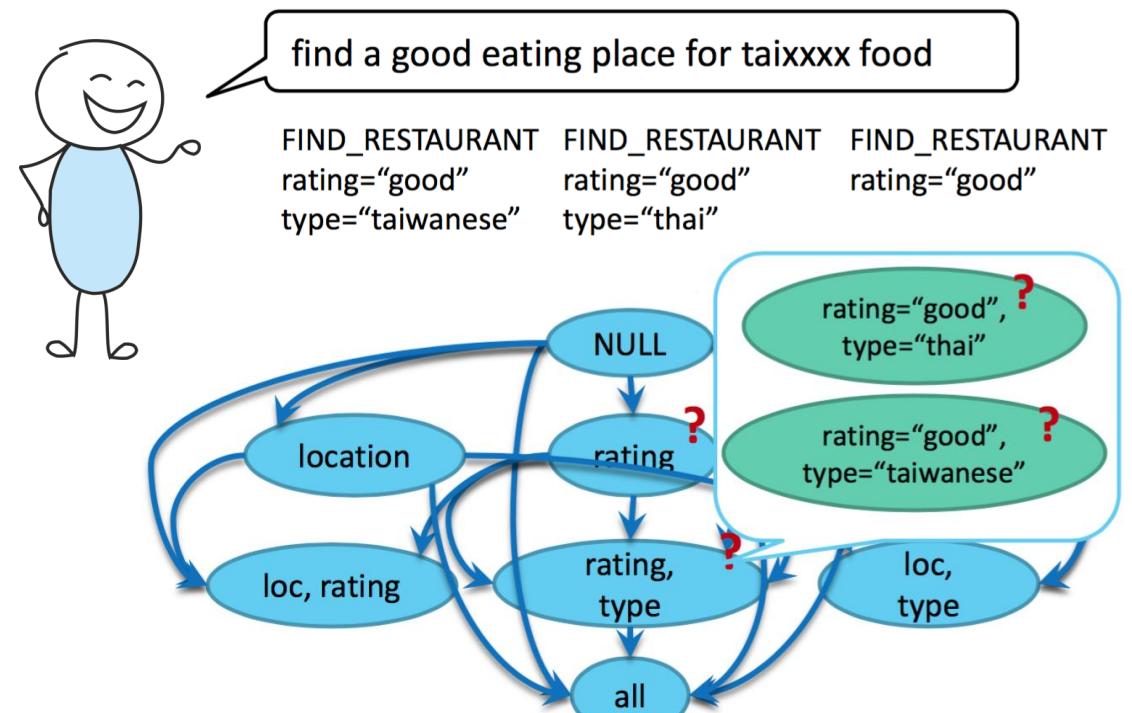
a: Cocum is a cheap
place in the West.
*inform(name=cocum, 0.3
area=west, 0.1
pricerange=cheap)*

u: What is their
number and
address?
*request(address,
phone)*

Dialogue State Tracking

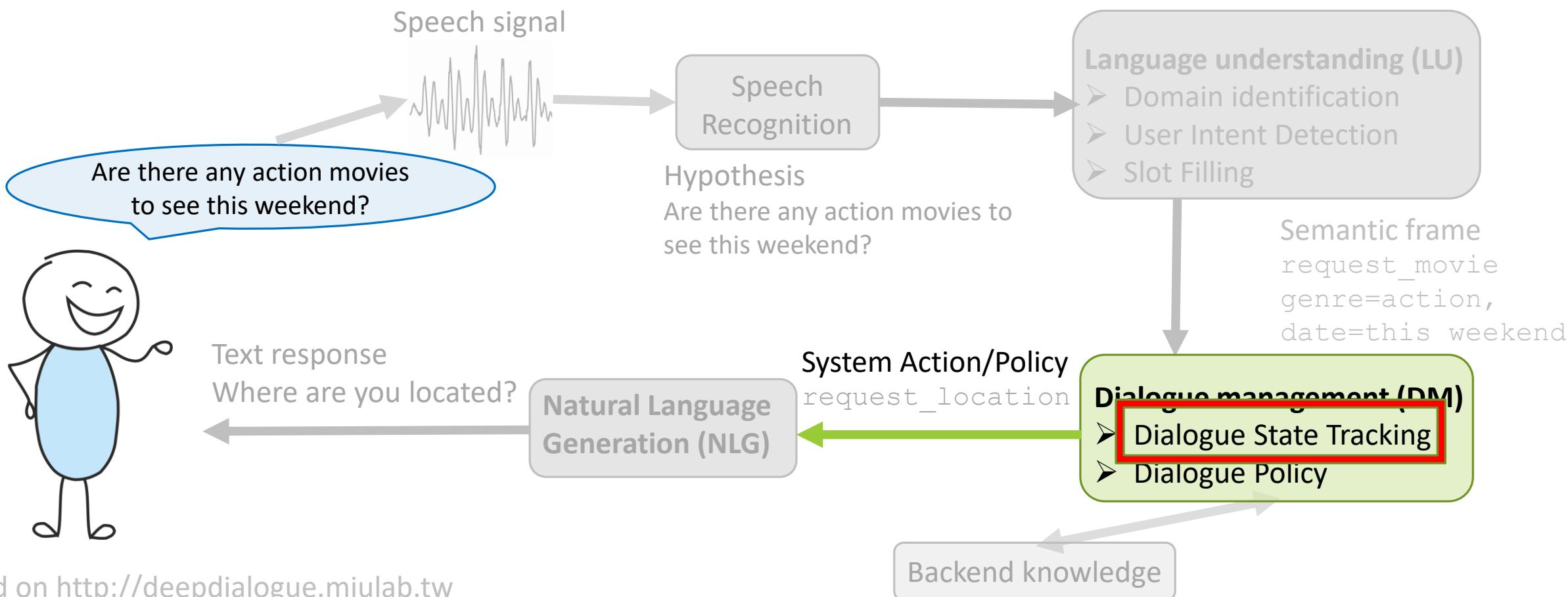
Maintain a distribution over multiple hypotheses of the true dialog state:

- robust sets of handcrafted rules
- conditional random fields
- maximum entropy models
- web-style ranking
- Deep learning: RNNs, LSTMs, etc

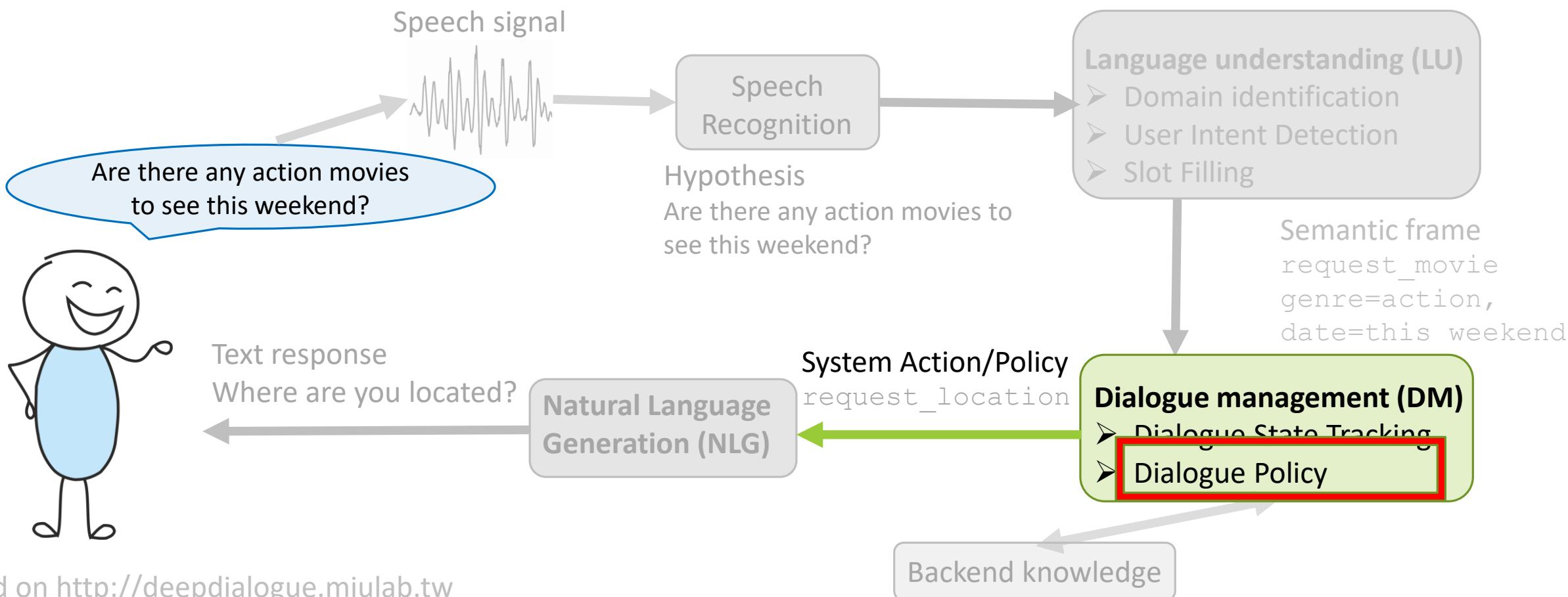


Picture from: <http://deepdialogue.miulab.tw>

Dialogue Management

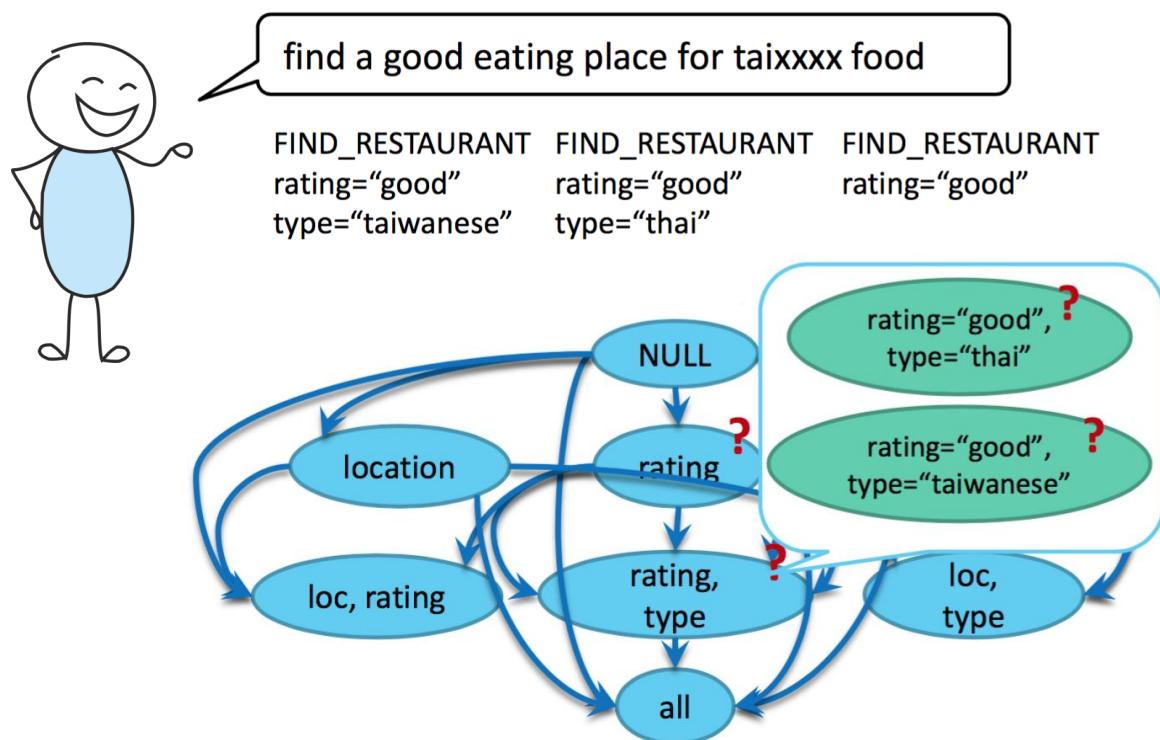


Dialogue Management



Dialogue Policy

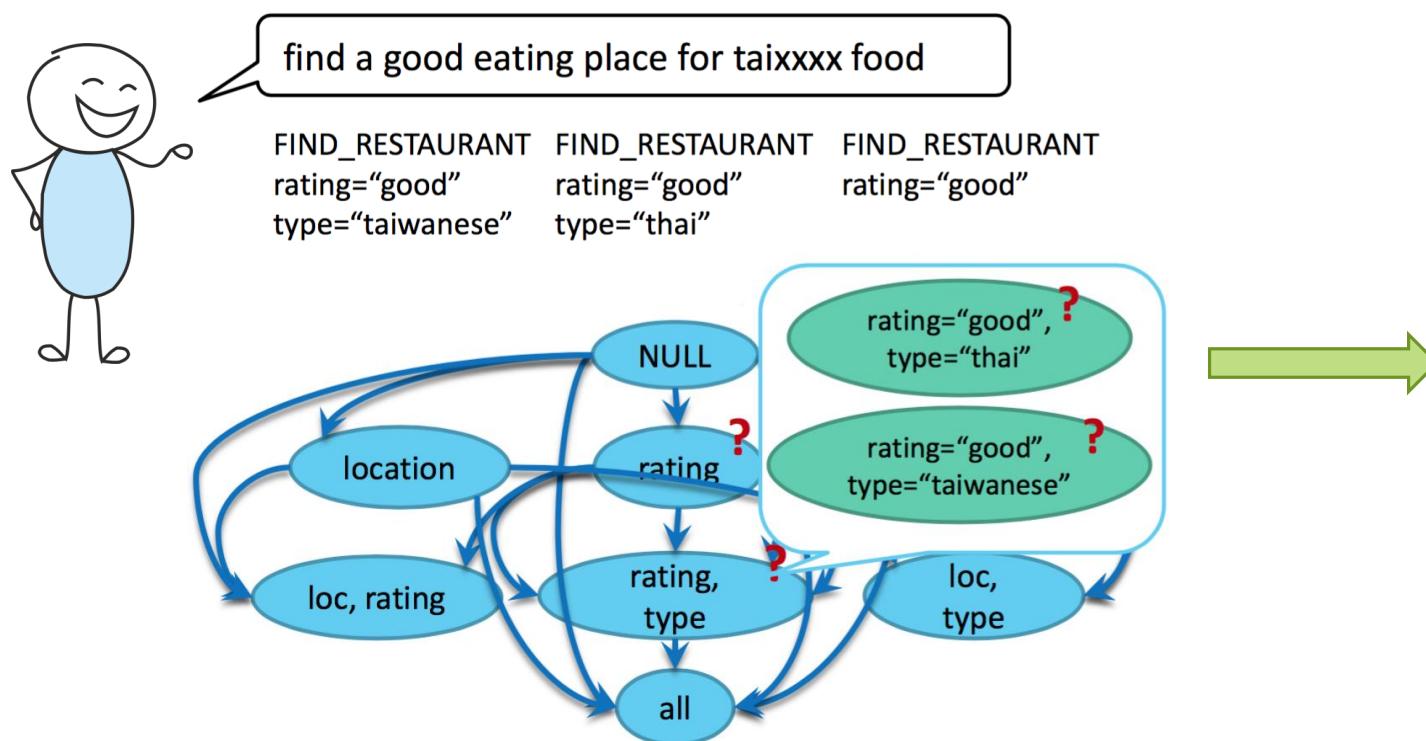
Dialogue State



Picture partially from: <http://deepdialogue.miulab.tw>

Dialogue Policy

Dialogue State

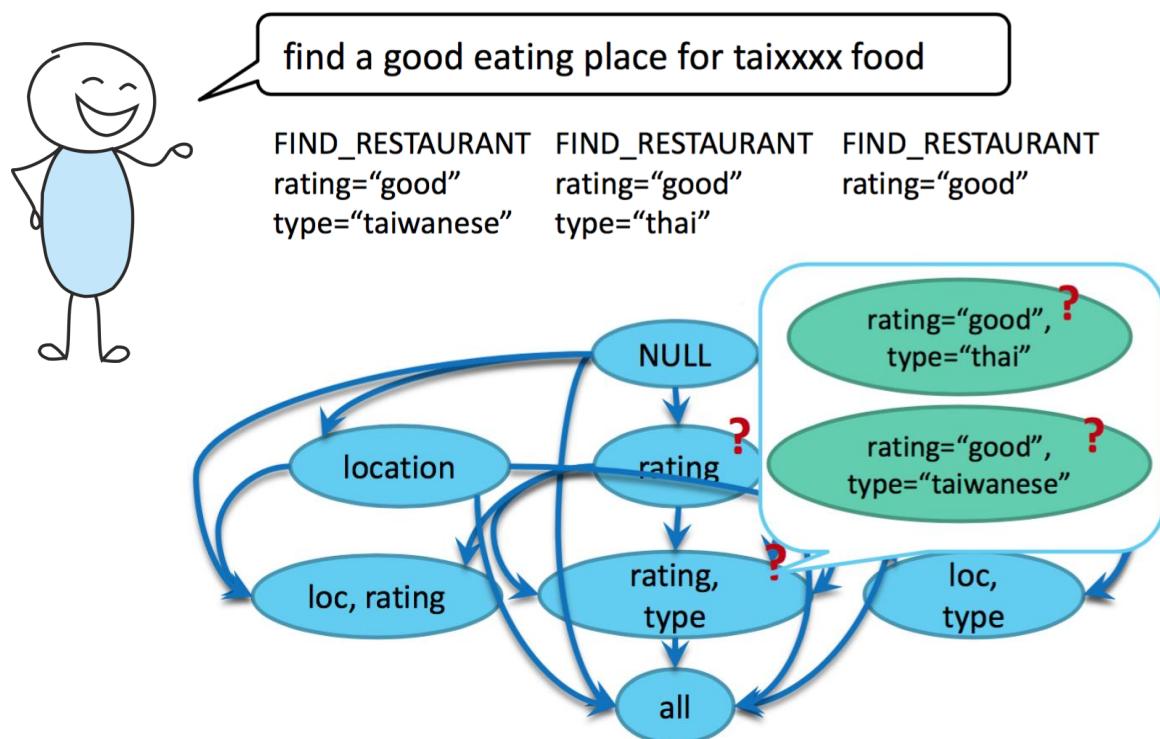


Agent Action

- Inform(location="Taipei 101")
 - "The nearest one is at Taipei 101"
- Request(location)
 - "Where is your home?"
- Confirm(type="taiwanese")
 - "Did you want Taiwanese food?"

Dialogue Policy

Dialogue State



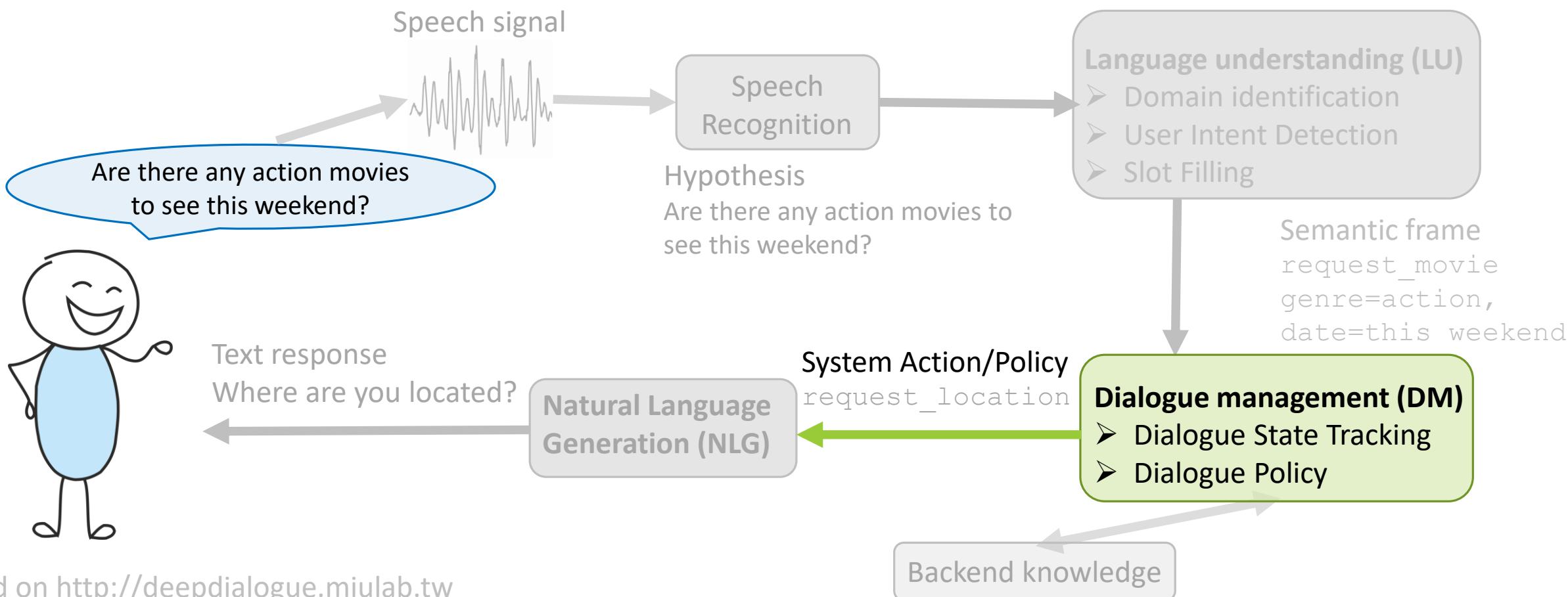
Agent Action

- Inform(location="Taipei 101")
 - "The nearest one is at Taipei 101"
- Request(location)
 - "Where is your home?"
- Confirm(type="taiwanese")
 - "Did you want Taiwanese food?"

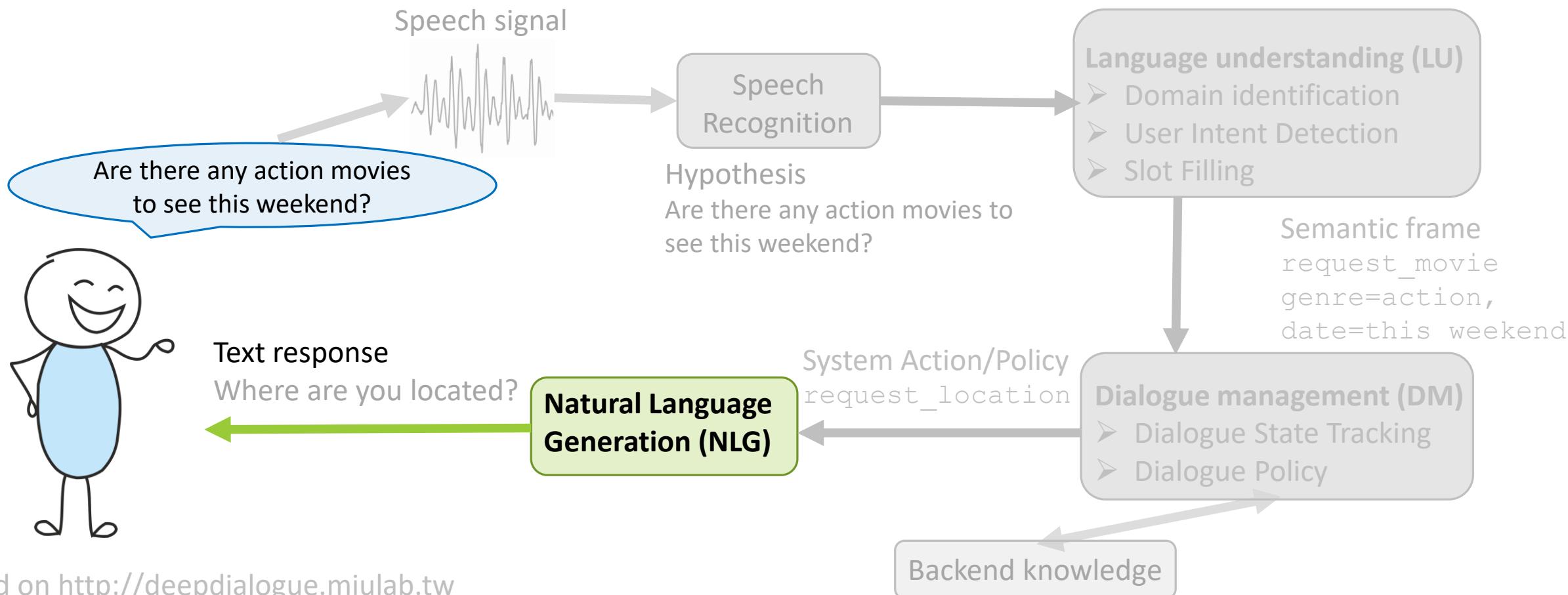
A lot of methods for optimizing policy will be in the RL course
(4 semester)!

Natural Language Generation

Dialogue Management



Natural Language Generation



Natural Language Generation

Generate natural language response given a dialogue action

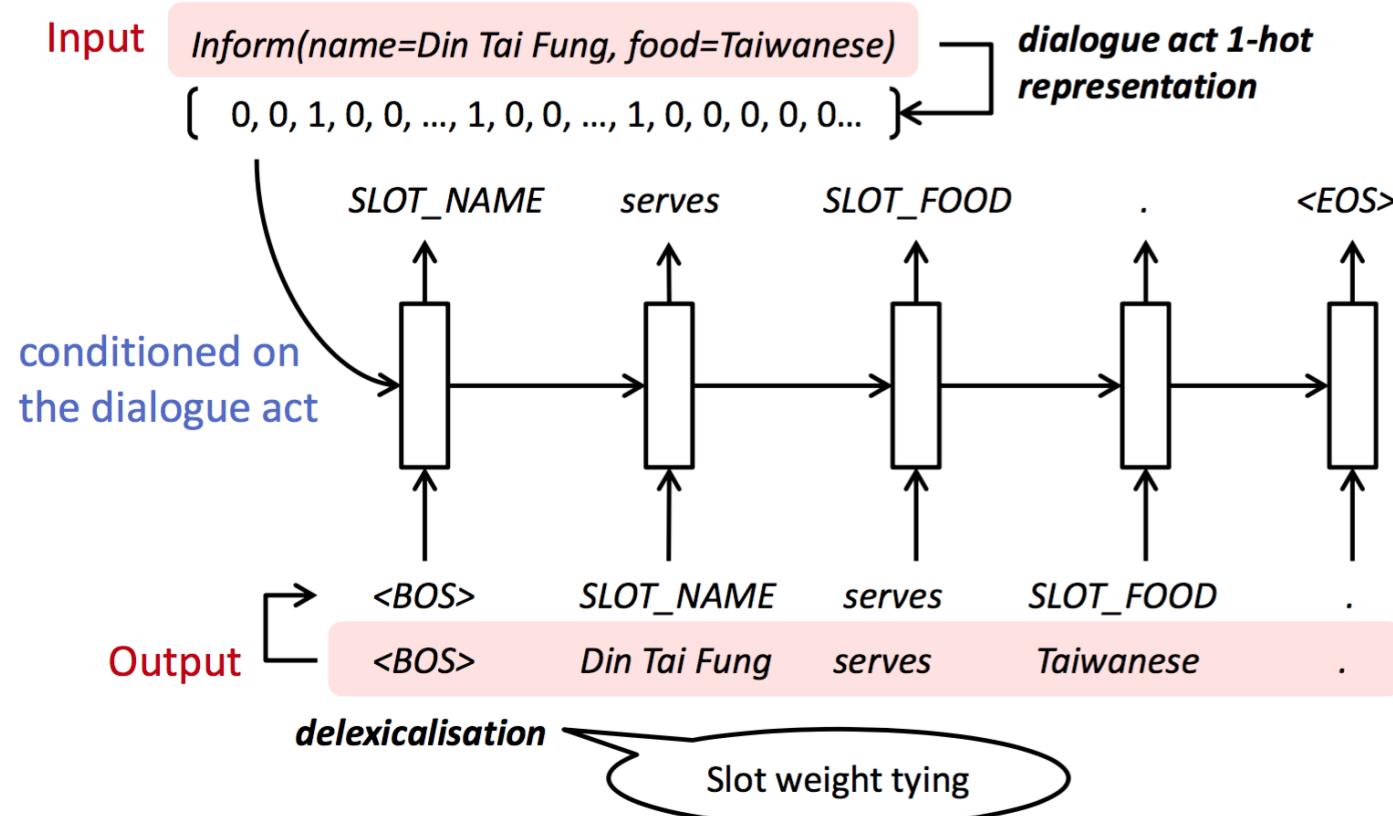
Inform(location="Taipei 101") → "The nearest one is at Taipei 101"

Request(location) → "Where is your home?"

Confirm(type="taiwanese") → "Did you want Taiwanese food?"

Example from: <http://deepdialogue.miulab.tw>

RNN-Based LM NLG

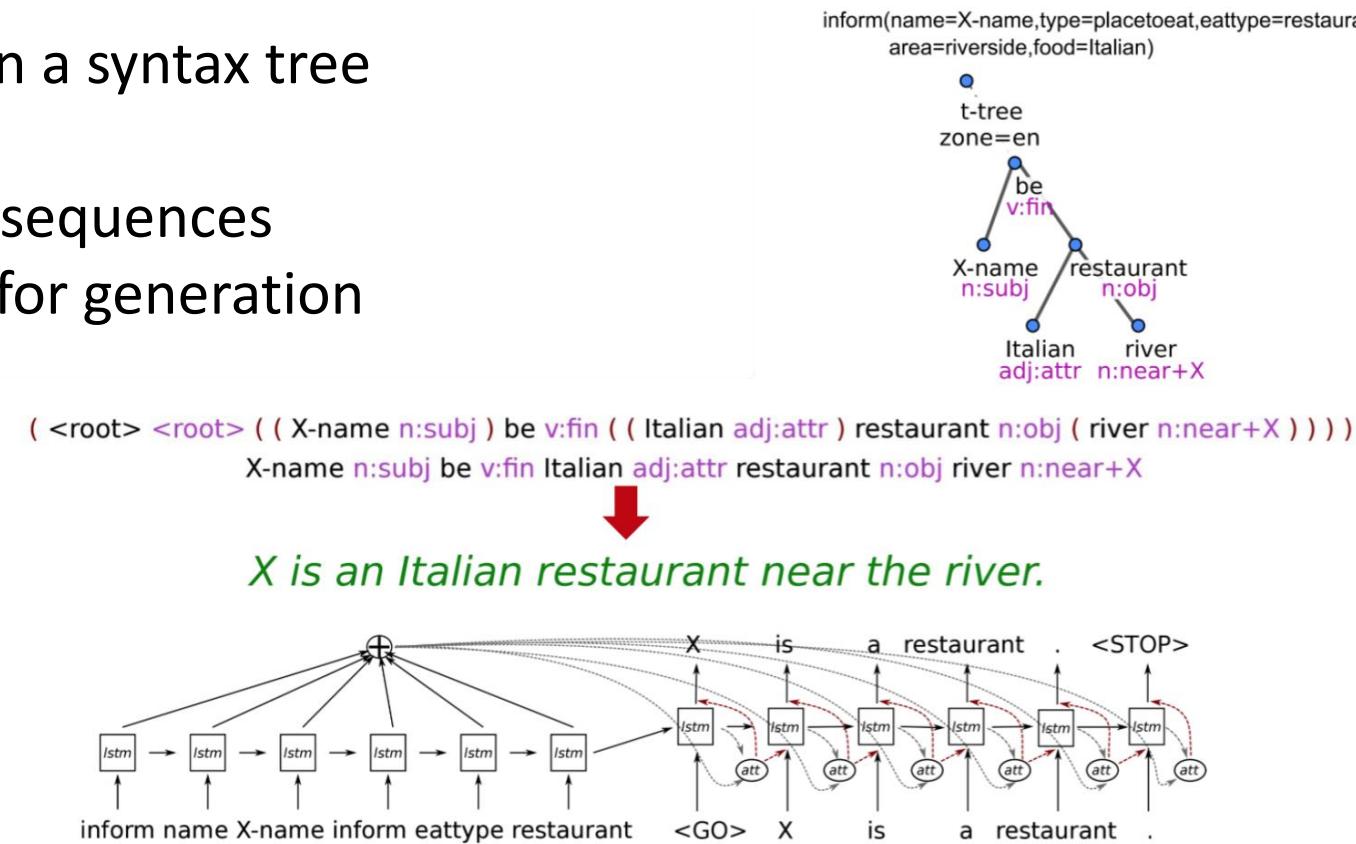


Picture from: <http://deepdialogue.miulab.tw> , paper: <http://www.aclweb.org/W/W15/W15-46.pdf#page=295>

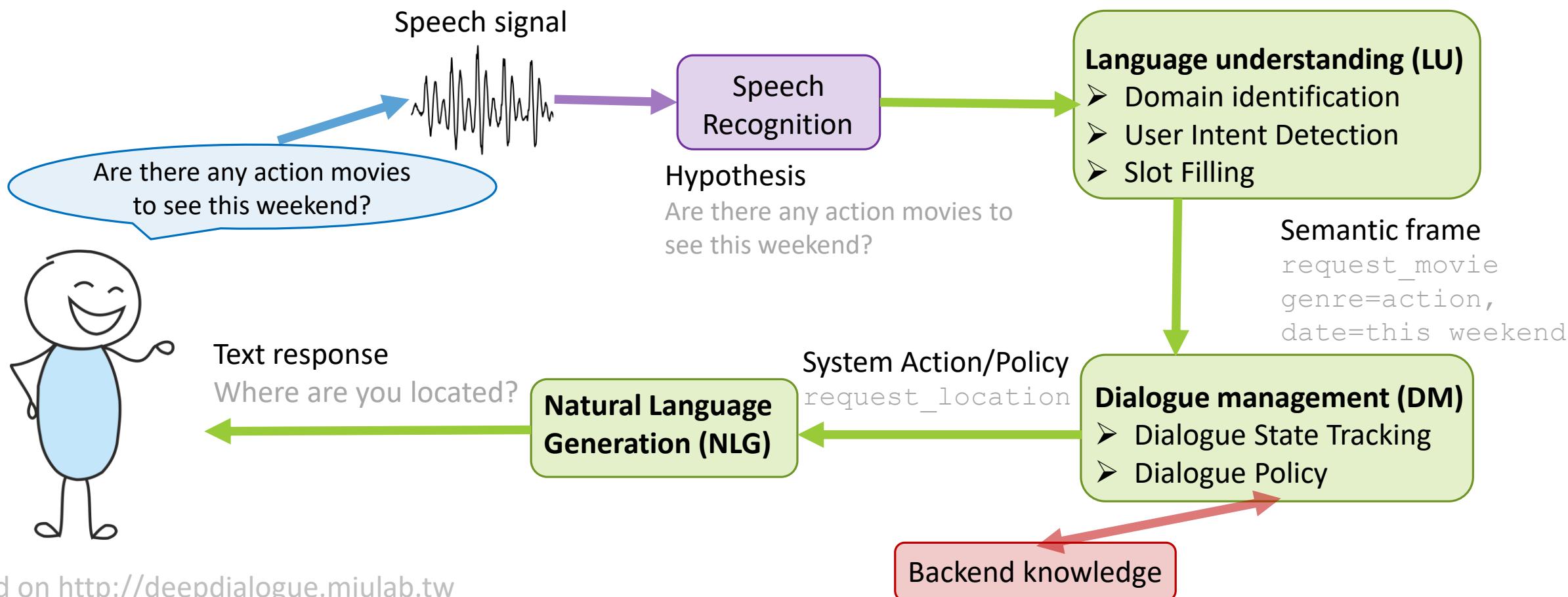
Structural NLG

Goal: NLG based on a syntax tree

- Encode trees as sequences
 - Seq2Seq model for generation



Task-oriented dialogue system



General conversation

AKA CHIT-CHATS

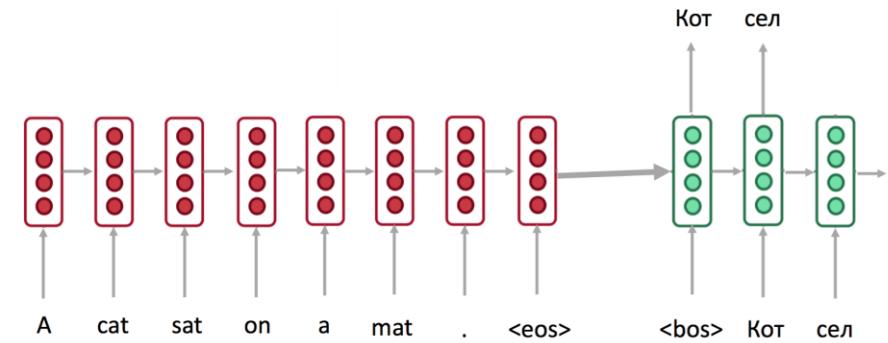
Retrieval-based (extractive)

Select a response from a database



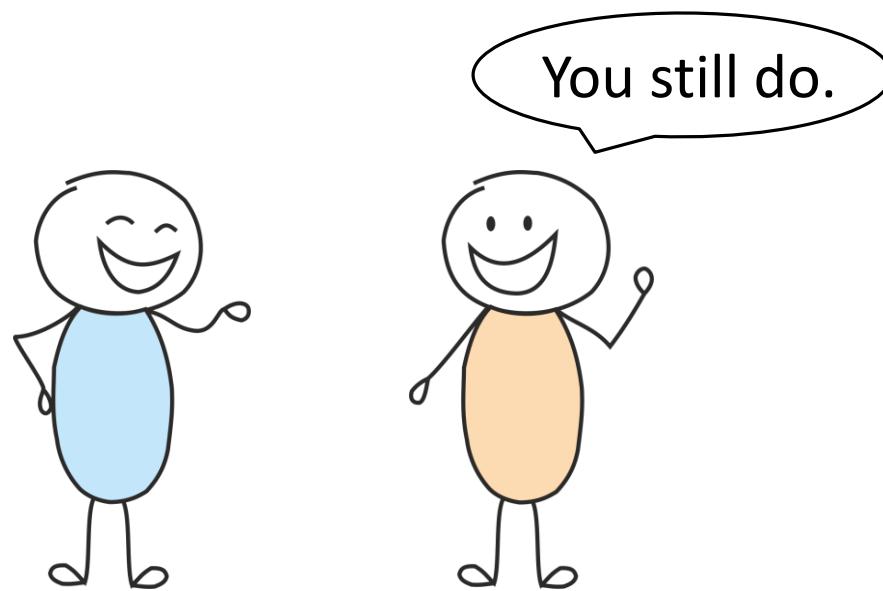
Generative (abstractive)

Generate a response
with some model

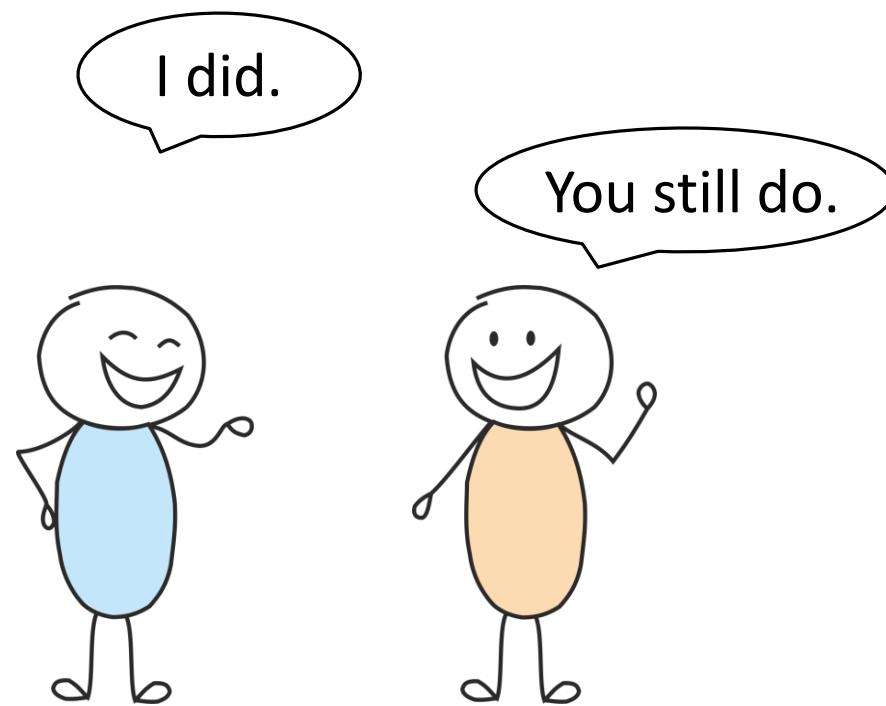


Generative Models for General conversation

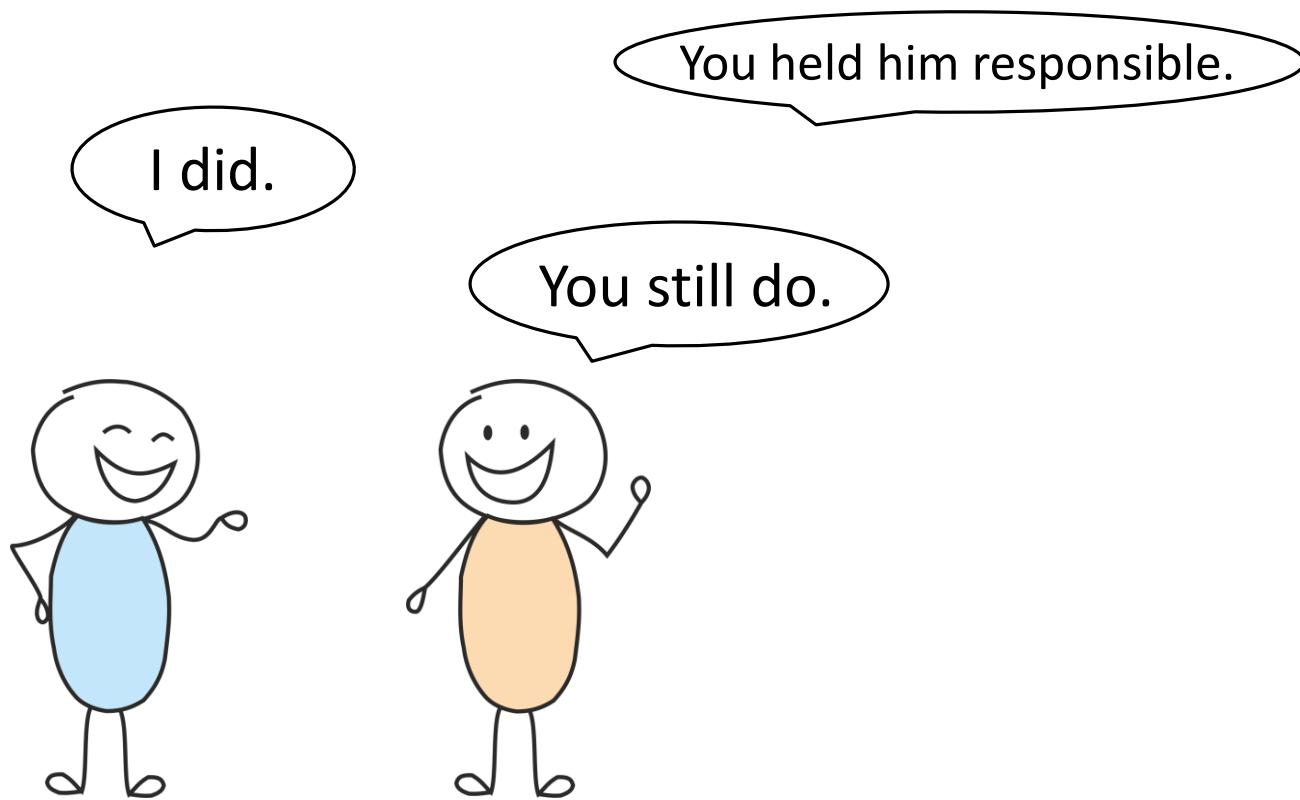
Storyline 1. The importance of context.



Storyline 1. The importance of context.



Storyline 1. The importance of context.



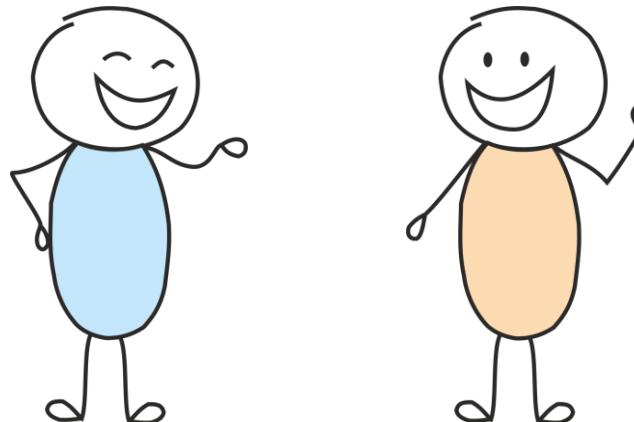
Storyline 1. The importance of context.

So I didn't speak to my father
for two years after that.

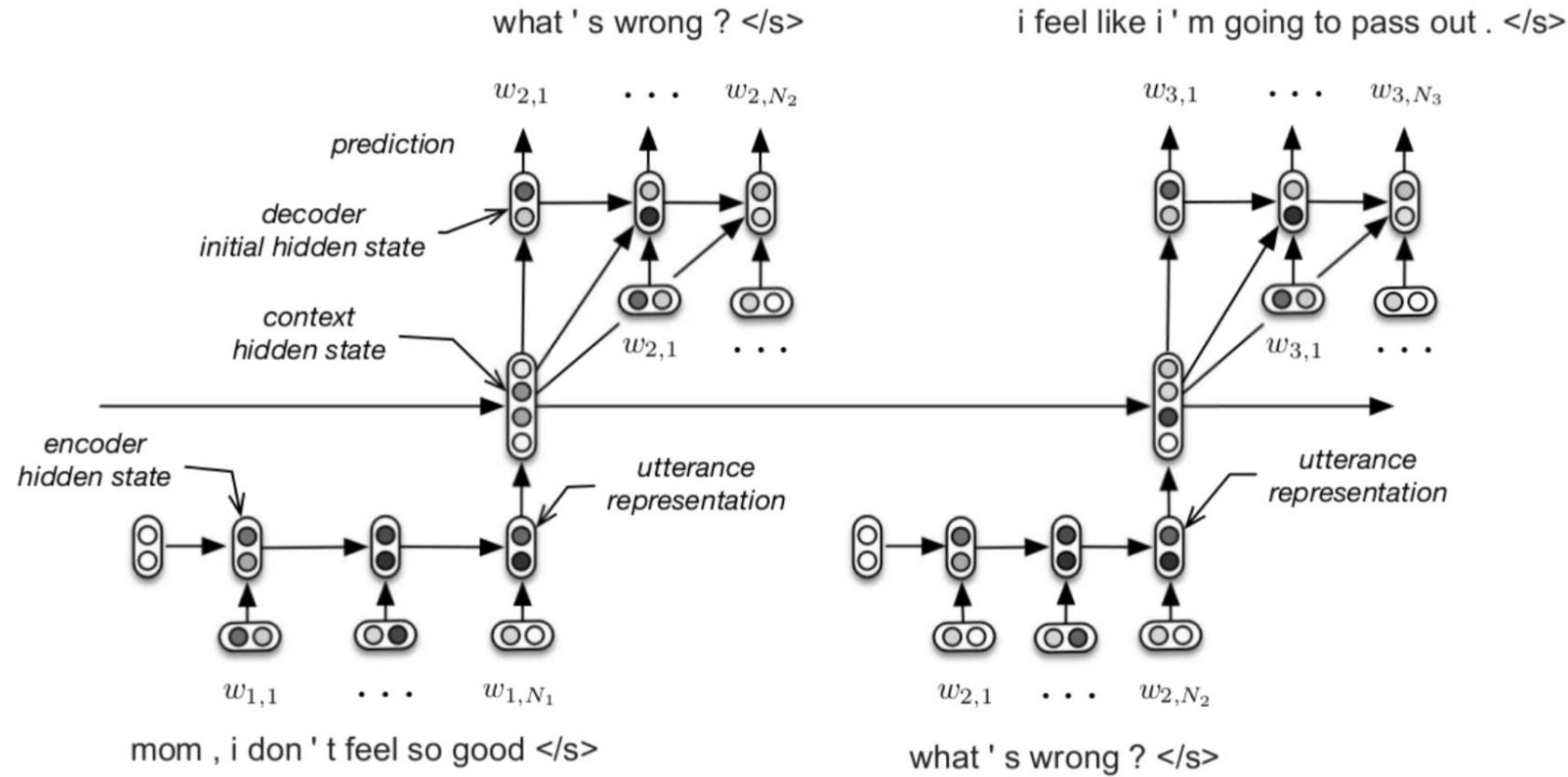
I did.

You held him responsible.

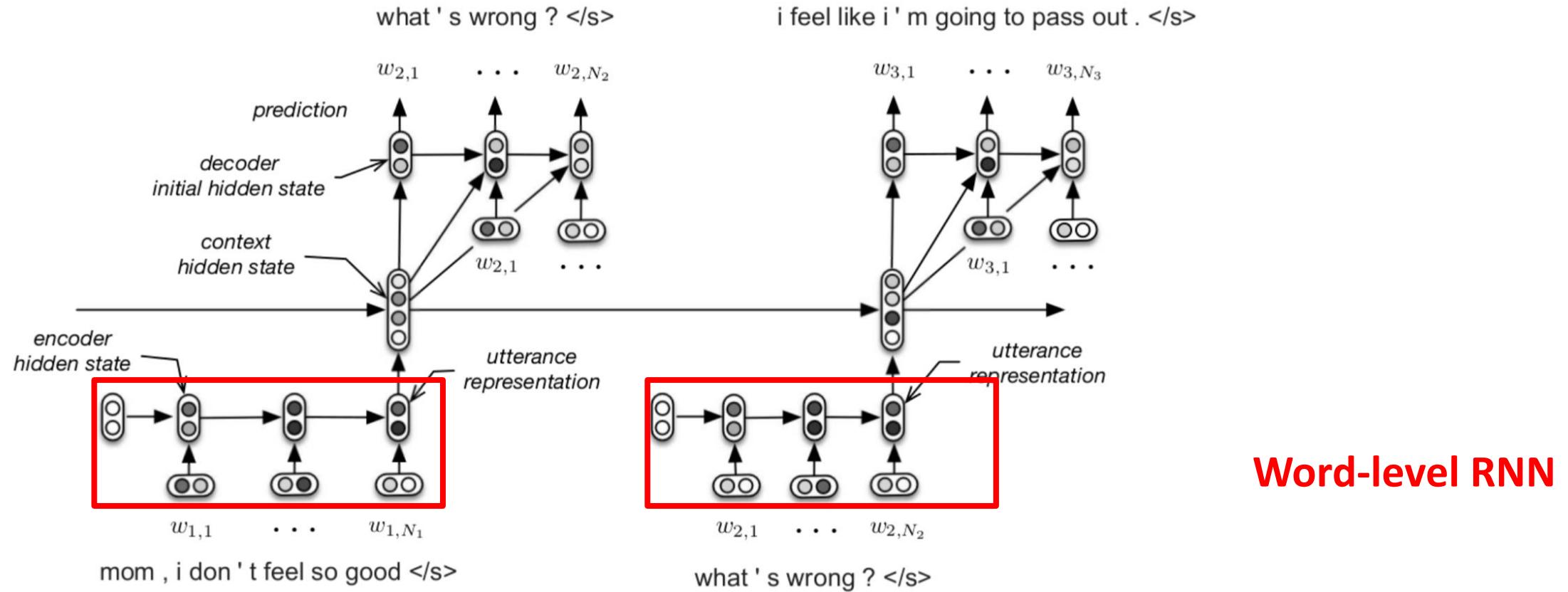
You still do.



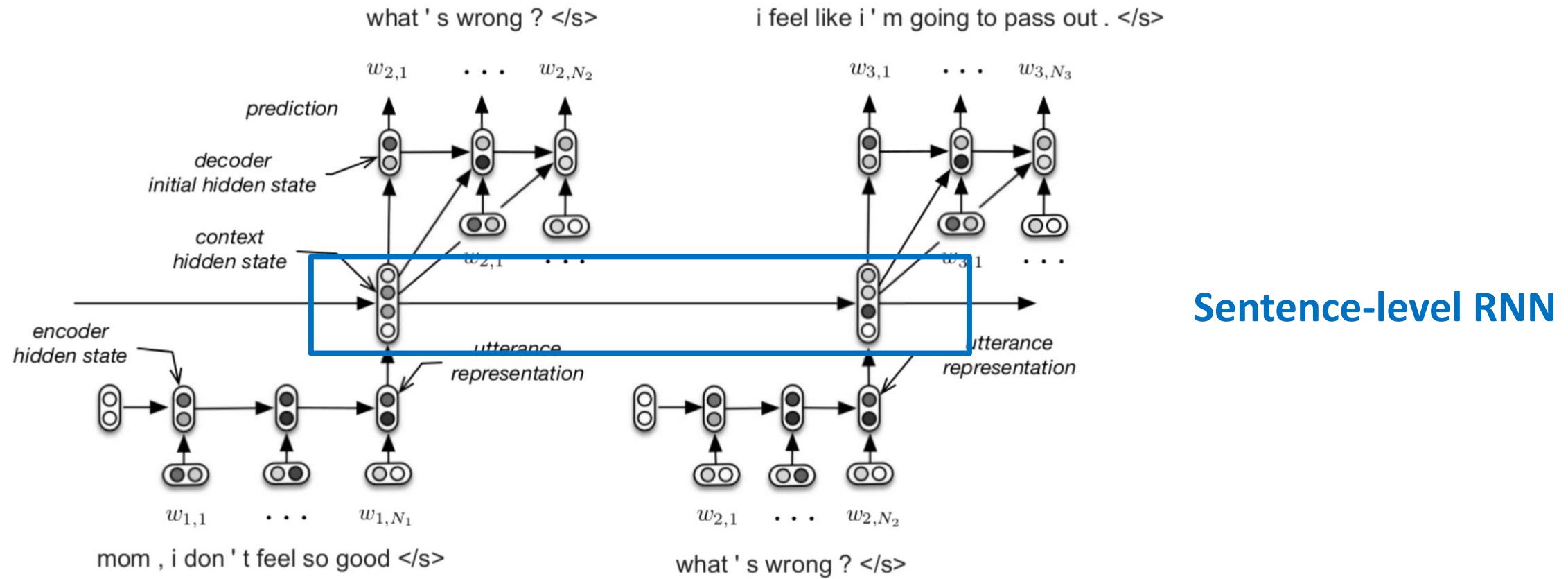
Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models



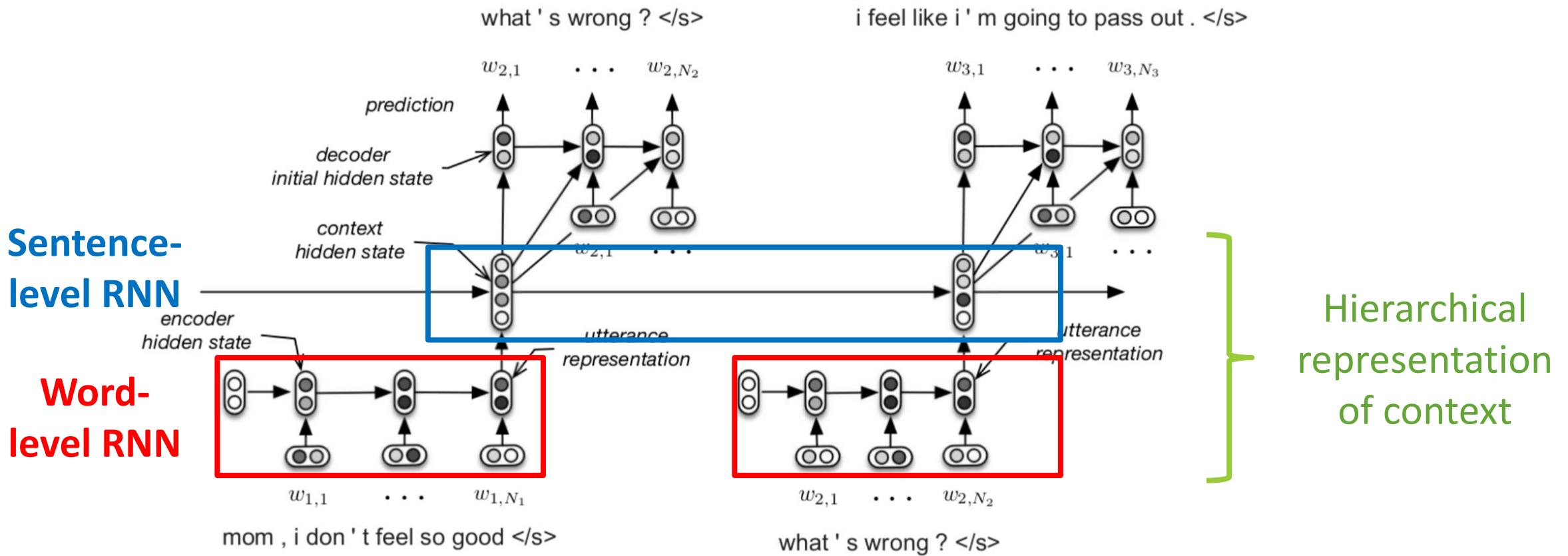
Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models



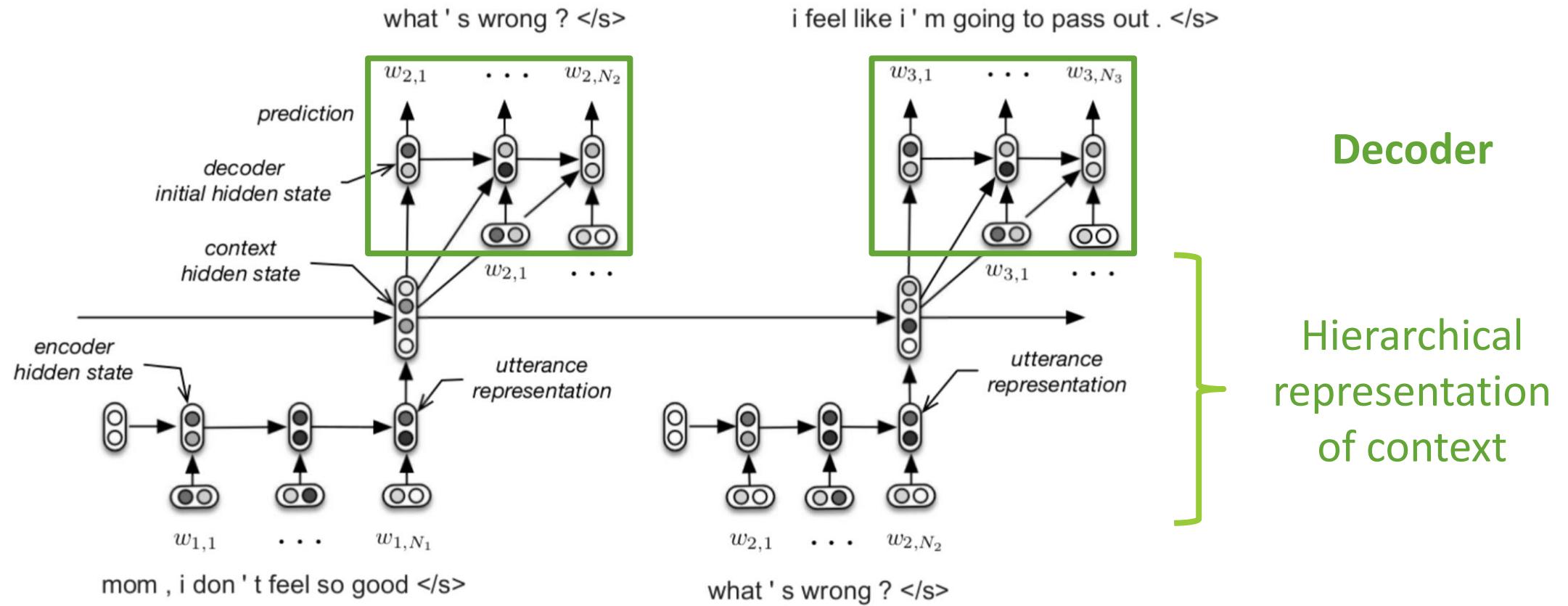
Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models



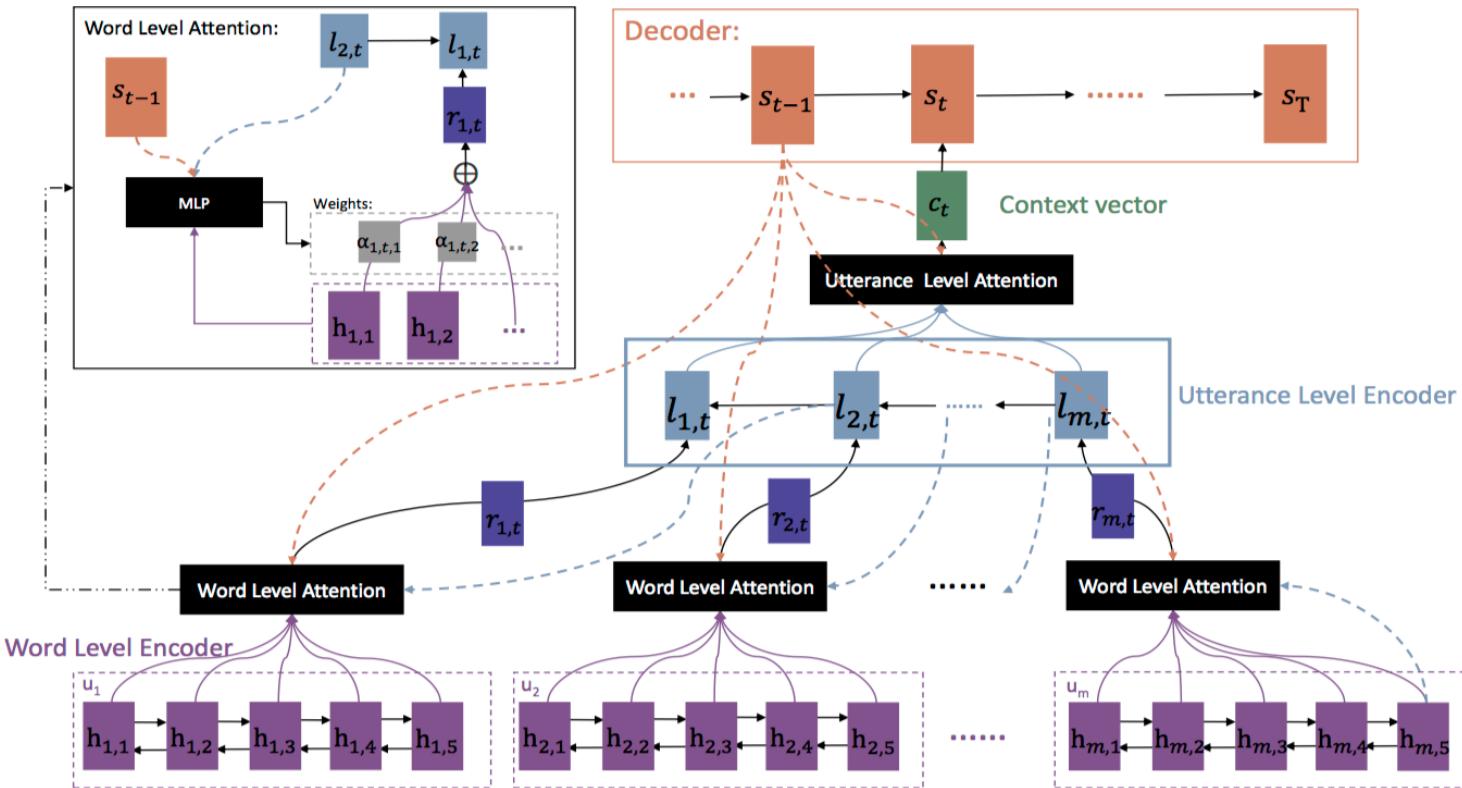
Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models



Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models

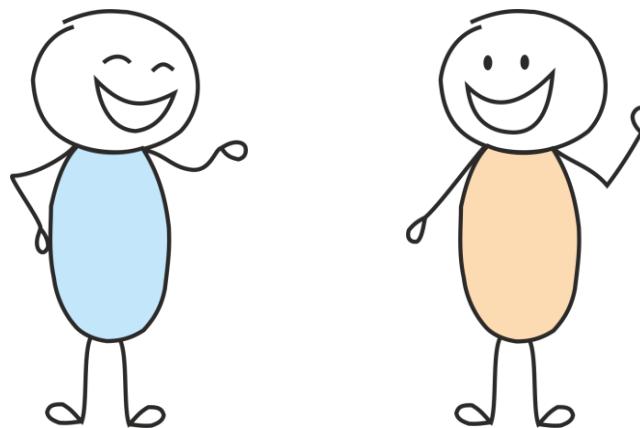


Hierarchical Recurrent Attention Network for Response Generation



Storyline 2.

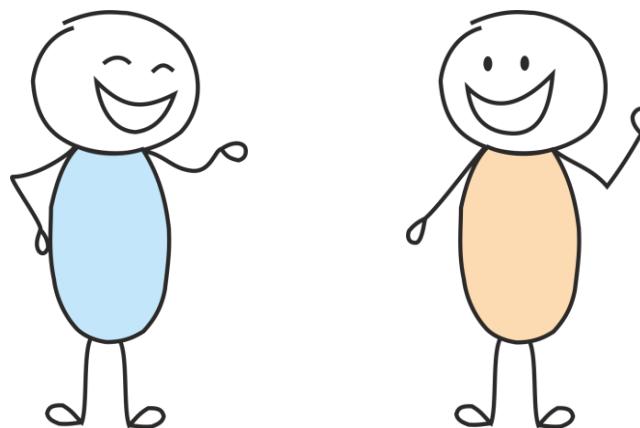
The importance of response diversity.



Storyline 2.

The importance of response diversity.

Hi! How are you doing?

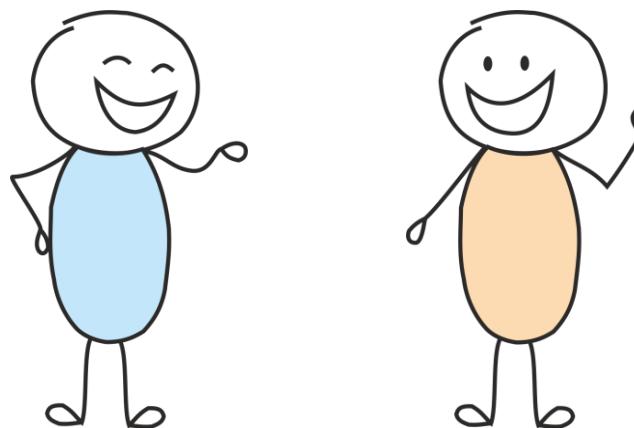


Storyline 2.

The importance of response diversity.

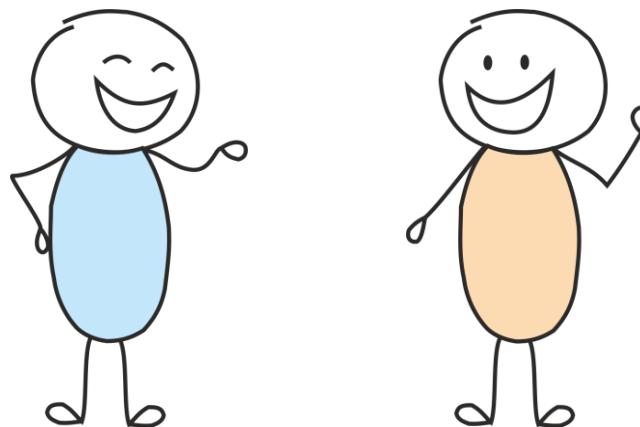
Hi! How are you doing?

I'm ok.



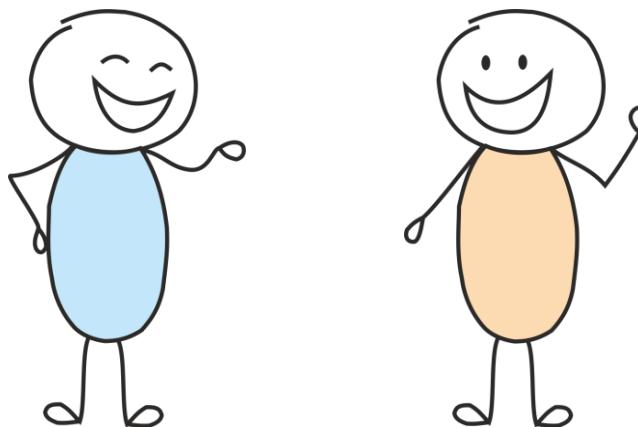
Storyline 2.

The importance of response diversity.



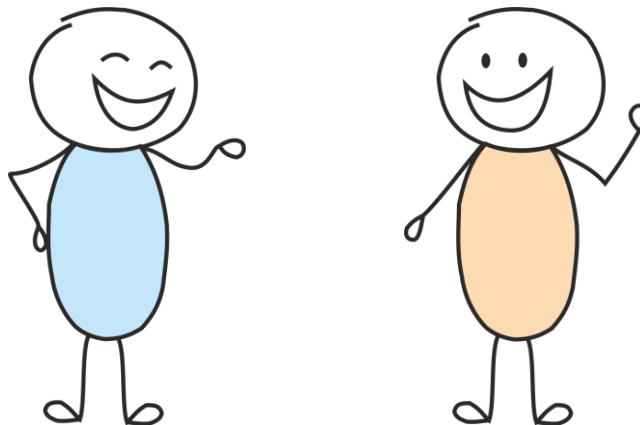
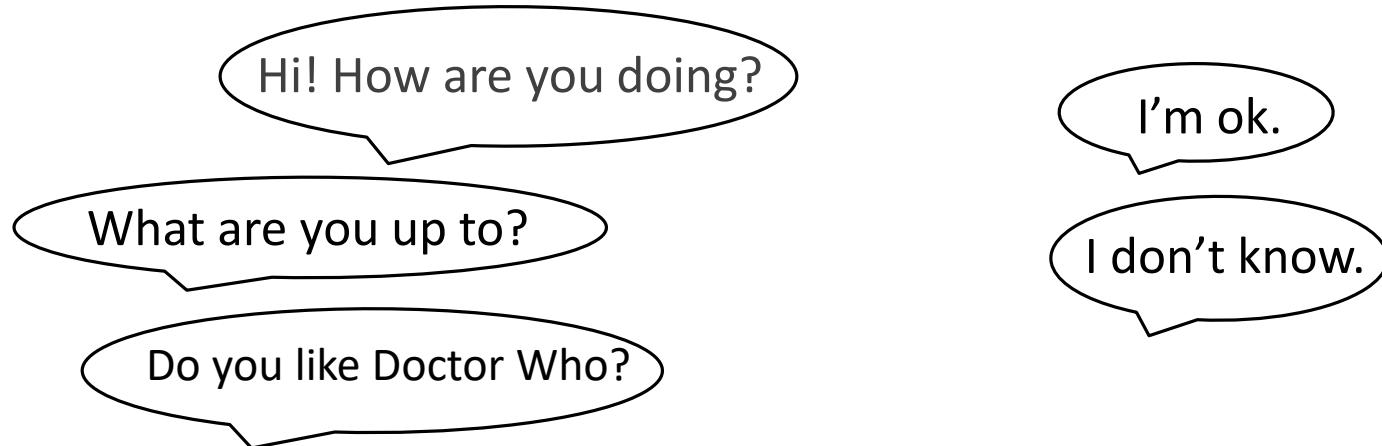
Storyline 2.

The importance of response diversity.



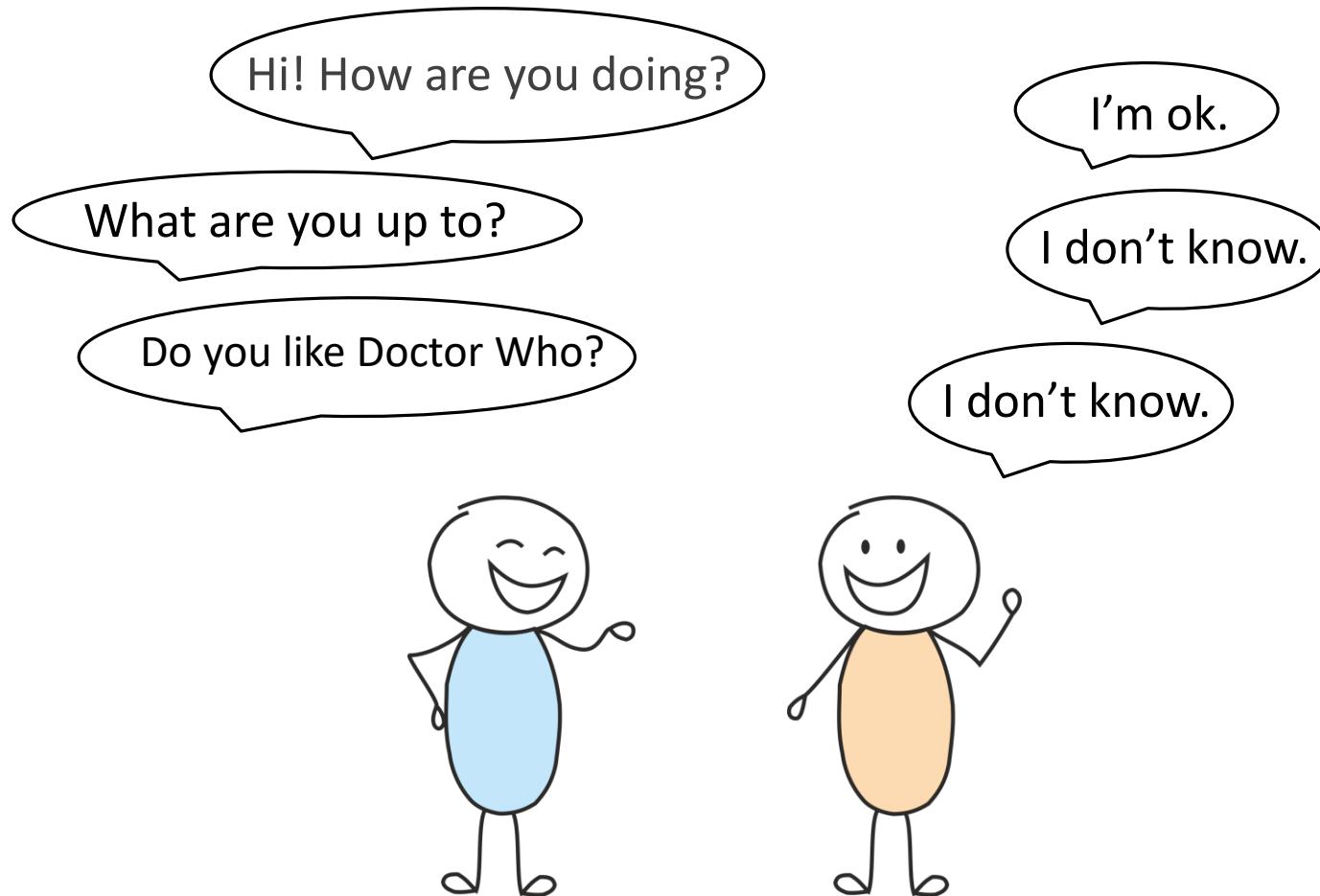
Storyline 2.

The importance of response diversity.



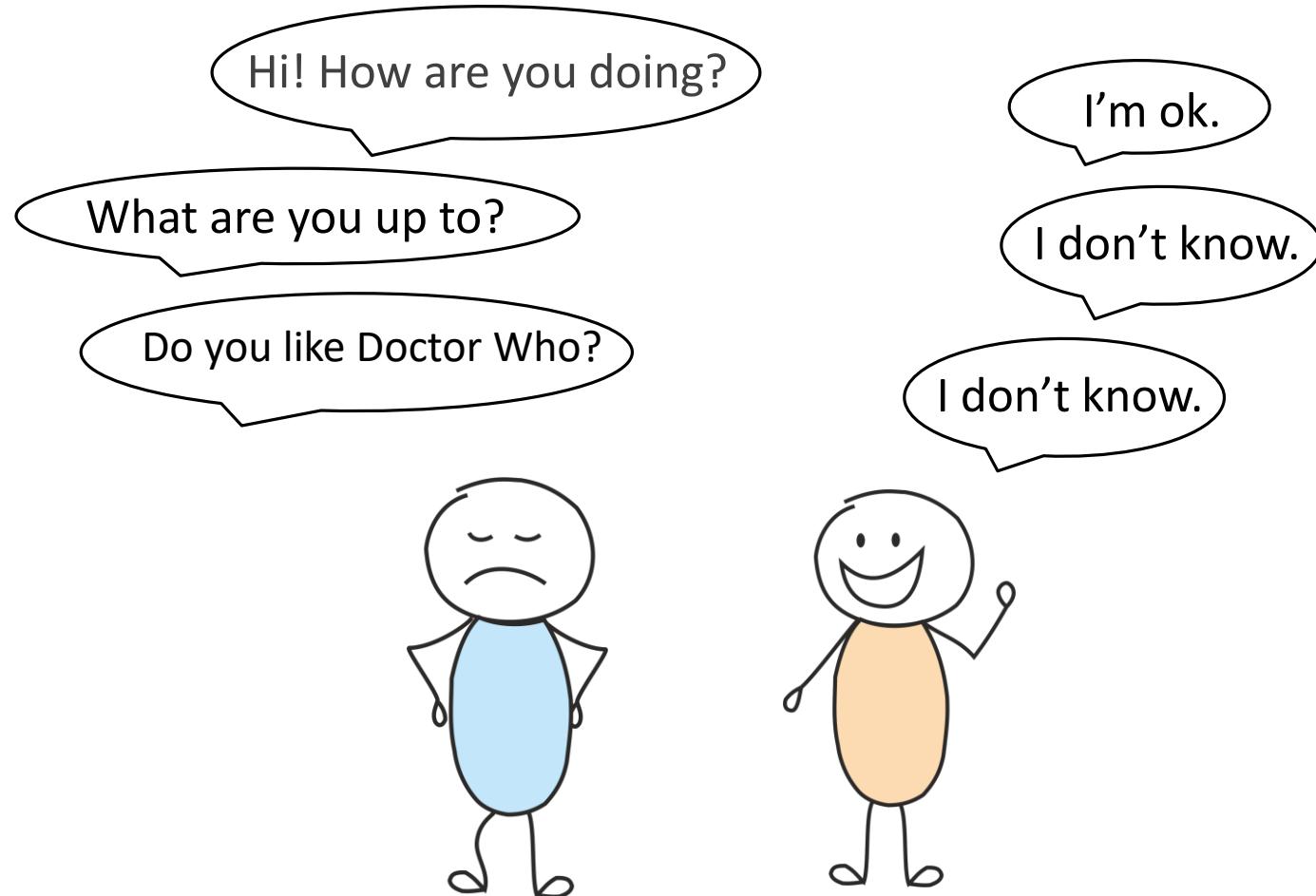
Storyline 2.

The importance of response diversity.



Storyline 2.

The importance of response diversity.



A Diversity-Promoting Objective Function for Neural Conversation Models

Input: What are you doing?

| | | | |
|-------|---------------------|-------|---------------------|
| -0.86 | I don't know. | -1.09 | Get out of here. |
| -1.03 | I don't know! | -1.09 | I'm going home. |
| -1.06 | Nothing. | -1.09 | Oh my god! |
| -1.09 | Get out of the way. | -1.10 | I'm talking to you. |

Input: what is your name?

| | | | |
|-------|--------------------|-------|--------------------|
| -0.91 | I don't know. | ... | |
| -0.92 | I don't know! | -1.55 | My name is Robert. |
| -0.92 | I don't know, sir. | -1.58 | My name is John. |
| -0.97 | Oh, my god! | -1.59 | My name's John. |

Input: How old are you?

| | | | |
|-------|----------------|-------|--------------|
| -0.79 | I don't know. | ... | |
| -1.06 | I'm fine. | -1.64 | Twenty-five. |
| -1.17 | I'm all right. | -1.66 | Five. |
| -1.17 | I'm not sure. | -1.71 | Eight. |

Table 1: Responses generated by a 4-layer SEQ2SEQ neural model trained on 20 million conversation pairs take from the OpenSubtitles dataset. Decoding is implemented with beam size set to 200. The top examples are the responses with the highest average probability log-likelihoods in the N-best list. Lower-ranked, less-generic responses were manually chosen.

Standard seq2seq: $\hat{T} = \arg \max_T \{ \log p(T|S) \}$

A Diversity-Promoting Objective Function for Neural Conversation Models

Input: What are you doing?

| | | | |
|-------|---------------------|-------|---------------------|
| -0.86 | I don't know. | -1.09 | Get out of here. |
| -1.03 | I don't know! | -1.09 | I'm going home. |
| -1.06 | Nothing. | -1.09 | Oh my god! |
| -1.09 | Get out of the way. | -1.10 | I'm talking to you. |

Input: what is your name?

| | | | |
|-------|--------------------|-------|--------------------|
| -0.91 | I don't know. | ... | |
| -0.92 | I don't know! | -1.55 | My name is Robert. |
| -0.92 | I don't know, sir. | -1.58 | My name is John. |
| -0.97 | Oh, my god! | -1.59 | My name's John. |

Input: How old are you?

| | | | |
|-------|----------------|-------|--------------|
| -0.79 | I don't know. | ... | |
| -1.06 | I'm fine. | -1.64 | Twenty-five. |
| -1.17 | I'm all right. | -1.66 | Five. |
| -1.17 | I'm not sure. | -1.71 | Eight. |

Table 1: Responses generated by a 4-layer SEQ2SEQ neural model trained on 20 million conversation pairs take from the OpenSubtitles dataset. Decoding is implemented with beam size set to 200. The top examples are the responses with the highest average probability log-likelihoods in the N-best list. Lower-ranked, less-generic responses were manually chosen.

Standard seq2seq: $\hat{T} = \arg \max_T \{ \log p(T|S) \}$

Mutual information
between source and target:

$$\log \frac{p(S, T)}{p(S)p(T)}$$

A Diversity-Promoting Objective Function for Neural Conversation Models

Input: What are you doing?

| | | | |
|-------|---------------------|-------|---------------------|
| -0.86 | I don't know. | -1.09 | Get out of here. |
| -1.03 | I don't know! | -1.09 | I'm going home. |
| -1.06 | Nothing. | -1.09 | Oh my god! |
| -1.09 | Get out of the way. | -1.10 | I'm talking to you. |

Input: what is your name?

| | | |
|-------|--------------------|--------------------------|
| -0.91 | I don't know. | ... |
| -0.92 | I don't know! | -1.55 My name is Robert. |
| -0.92 | I don't know, sir. | -1.58 My name is John. |
| -0.97 | Oh, my god! | -1.59 My name's John. |

Input: How old are you?

| | | |
|-------|----------------|--------------------|
| -0.79 | I don't know. | ... |
| -1.06 | I'm fine. | -1.64 Twenty-five. |
| -1.17 | I'm all right. | -1.66 Five. |
| -1.17 | I'm not sure. | -1.71 Eight. |

Table 1: Responses generated by a 4-layer SEQ2SEQ neural model trained on 20 million conversation pairs take from the OpenSubtitles dataset. Decoding is implemented with beam size set to 200. The top examples are the responses with the highest average probability log-likelihoods in the N-best list. Lower-ranked, less-generic responses were manually chosen.

Standard seq2seq: $\hat{T} = \arg \max_T \{ \log p(T|S) \}$

Mutual information
between source and target:

$$\log \frac{p(S, T)}{p(S)p(T)}$$

Maximum Mutual
Information (MMI): $\hat{T} = \arg \max_T \{ \log p(T|S) - \log p(T) \}$

A Diversity-Promoting Objective Function for Neural Conversation Models

Input: What are you doing?

| | | | |
|-------|---------------------|-------|---------------------|
| -0.86 | I don't know. | -1.09 | Get out of here. |
| -1.03 | I don't know! | -1.09 | I'm going home. |
| -1.06 | Nothing. | -1.09 | Oh my god! |
| -1.09 | Get out of the way. | -1.10 | I'm talking to you. |

Input: what is your name?

| | | |
|-------|--------------------|--------------------------|
| -0.91 | I don't know. | ... |
| -0.92 | I don't know! | -1.55 My name is Robert. |
| -0.92 | I don't know, sir. | -1.58 My name is John. |
| -0.97 | Oh, my god! | -1.59 My name's John. |

Input: How old are you?

| | | |
|-------|----------------|--------------------|
| -0.79 | I don't know. | ... |
| -1.06 | I'm fine. | -1.64 Twenty-five. |
| -1.17 | I'm all right. | -1.66 Five. |
| -1.17 | I'm not sure. | -1.71 Eight. |

Table 1: Responses generated by a 4-layer SEQ2SEQ neural model trained on 20 million conversation pairs take from the OpenSubtitles dataset. Decoding is implemented with beam size set to 200. The top examples are the responses with the highest average probability log-likelihoods in the N-best list. Lower-ranked, less-generic responses were manually chosen.

Standard seq2seq: $\hat{T} = \arg \max_T \{ \log p(T|S) \}$

Mutual information
between source and target:

$$\log \frac{p(S, T)}{p(S)p(T)}$$

Maximum Mutual
Information (MMI):

$$\hat{T} = \arg \max_T \{ \log p(T|S) - \log p(T) \}$$

Penalize frequent answers

A Diversity-Promoting Objective Function for Neural Conversation Models

Input: What are you doing?

| | | | |
|-------|---------------------|-------|---------------------|
| -0.86 | I don't know. | -1.09 | Get out of here. |
| -1.03 | I don't know! | -1.09 | I'm going home. |
| -1.06 | Nothing. | -1.09 | Oh my god! |
| -1.09 | Get out of the way. | -1.10 | I'm talking to you. |

Input: what is your name?

| | | |
|-------|--------------------|--------------------------|
| -0.91 | I don't know. | ... |
| -0.92 | I don't know! | -1.55 My name is Robert. |
| -0.92 | I don't know, sir. | -1.58 My name is John. |
| -0.97 | Oh, my god! | -1.59 My name's John. |

Input: How old are you?

| | | |
|-------|----------------|--------------------|
| -0.79 | I don't know. | ... |
| -1.06 | I'm fine. | -1.64 Twenty-five. |
| -1.17 | I'm all right. | -1.66 Five. |
| -1.17 | I'm not sure. | -1.71 Eight. |

Table 1: Responses generated by a 4-layer SEQ2SEQ neural model trained on 20 million conversation pairs take from the OpenSubtitles dataset. Decoding is implemented with beam size set to 200. The top examples are the responses with the highest average probability log-likelihoods in the N-best list. Lower-ranked, less-generic responses were manually chosen.

Standard seq2seq: $\hat{T} = \arg \max_T \{ \log p(T|S) \}$

Mutual information
between source and target:

$$\log \frac{p(S, T)}{p(S)p(T)}$$

Maximum Mutual
Information (MMI): $\hat{T} = \arg \max_T \{ \log p(T|S) - \log p(T) \}$

Generalized version: $\hat{T} = \arg \max_T \{ \log p(T|S) - \lambda \log p(T) \}$

Bayes Theorem: $\hat{T} = \arg \max_T \{ (1 - \lambda) \log p(T|S) + \lambda \log p(S|T) \}$

A Diversity-Promoting Objective Function for Neural Conversation Models

| message | SEQ2SEQ | $\text{MMI } P(T S) - \lambda P(T)$ |
|----------------------------------------------------------------------|--------------------------------------------|-----------------------------------------|
| I did not get the report from an mi-6 agent. | I don't think that's a good idea. | You did the right thing, did you? |
| You haven't been given an assignment in this case. | I don't know what you are talking about. | I've been looking all over for you. |
| I'm losing my grip. | I don't know what you are talking about. | I'm the only one in the world. |
| I am ready to help. | Come on, come on | I have something we need to talk about. |
| You programmed me to gather intelligence. That's all I've ever done. | You know that, don't you? | You do have fun, don't you? |
| I mean, we'd have to talk to him. | I mean, I don't know. | I mean, he's a good guy. |
| This has nothing to do with him. | You don't know anything about him, do you? | I know him as much. |
| How come you never say it? | I don't know | Because I don't want to hurt you |
| He is the one behind this scheme to crash x-net. | I don't know what you are talking about | Yes, he is, isn't he? |
| How much time do you have here? | I don't know. | Not long enough. Sorry, sir. |

Table 4: Responses from the SEQ2SEQ baseline and MMI-antiLM models on the OpenSubtitles dataset.

Diverse Beam Search

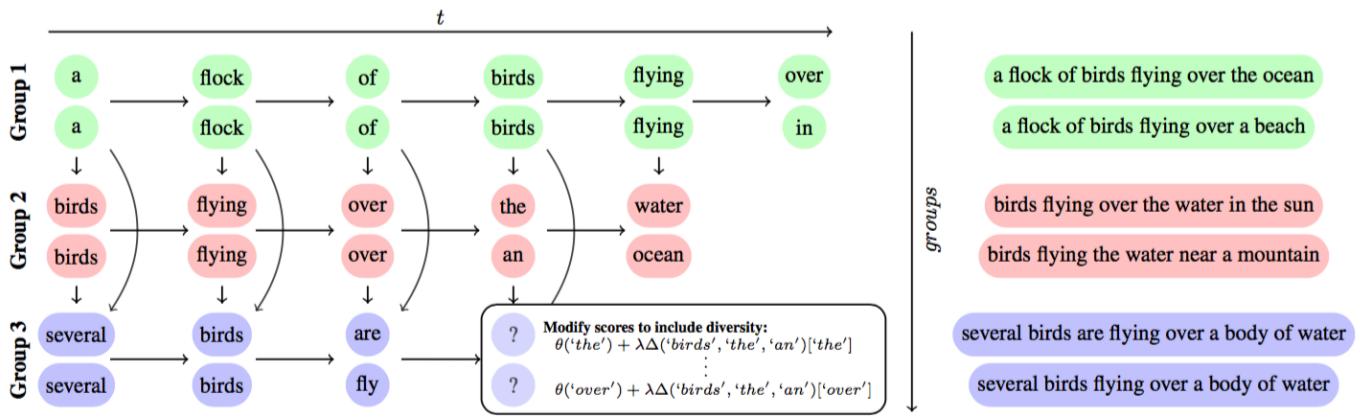


Figure 2: Diverse beam search operates left-to-right through time and top to bottom through groups. Diversity between groups is combined with joint log-probabilities, allowing continuations to be found efficiently. The resulting outputs are more diverse than for standard approaches.

- Divide beam on several groups
- Optimize each group keeping other fixed
- Add dissimilarity term

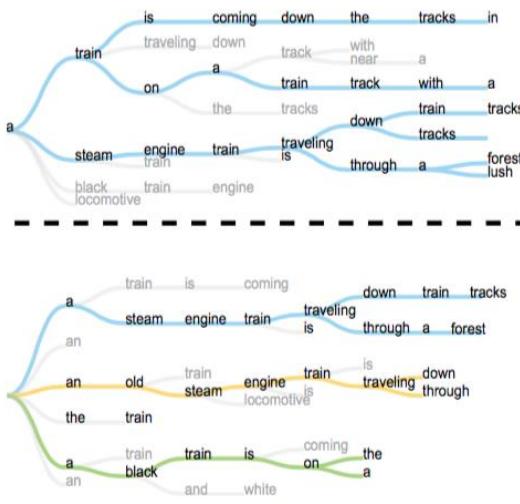
$$\Delta(\mathbf{y}_{[t]}, Y_{[t]}^g) = \sum_{b=1}^{B'} \delta(\mathbf{y}_{[t]}, \mathbf{y}_{b,[t]}^g)$$

Diverse Beam Search



Ground Truth Captions

Single engine train rolling down the tracks.
A steam locomotive is blowing steam.



Beam Search

A steam engine train travelling down train tracks.
A steam engine train travelling down tracks.
A steam engine train travelling through a forest.
A steam engine train travelling through a lush green forest.
A steam engine train travelling through a lush green countryside
A train on a train track with a sky background.

Diverse Beam Search

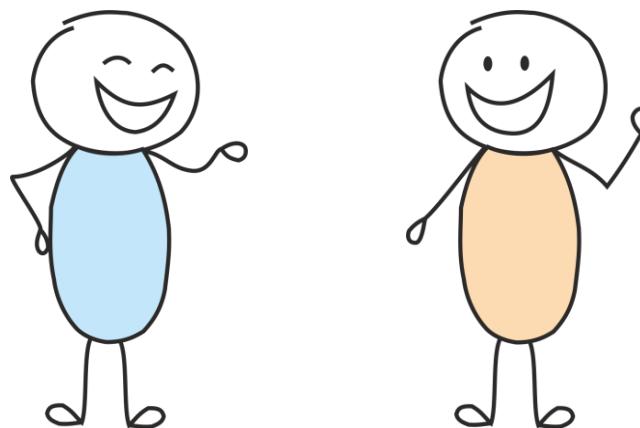
A steam engine travelling down train tracks.
A steam engine train travelling through a forest.
An old steam engine train travelling down train tracks.
An old steam engine train travelling through a forest.
A black train is on the tracks in a wooded area.
A black train is on the tracks in a rural area.

A locomotive drives along the tracks amongst trees and bushes.
An old fashion train with steam coming out of its pipe.
An engine is coming down the train track.
A black and red train moving down a train track.

Storyline 3.

The importance of coherence.

Are you a boy or a girl?

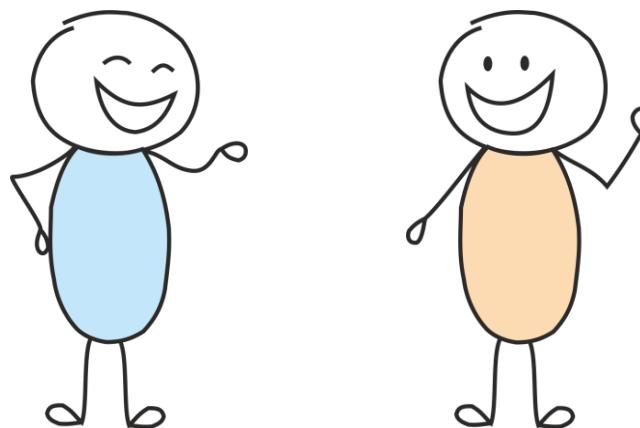


Storyline 3.

The importance of coherence.

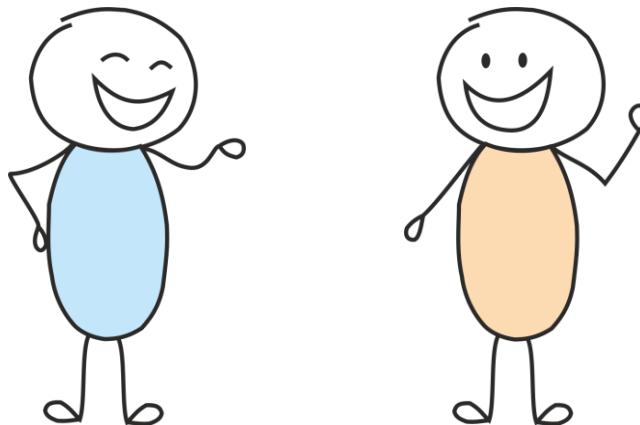
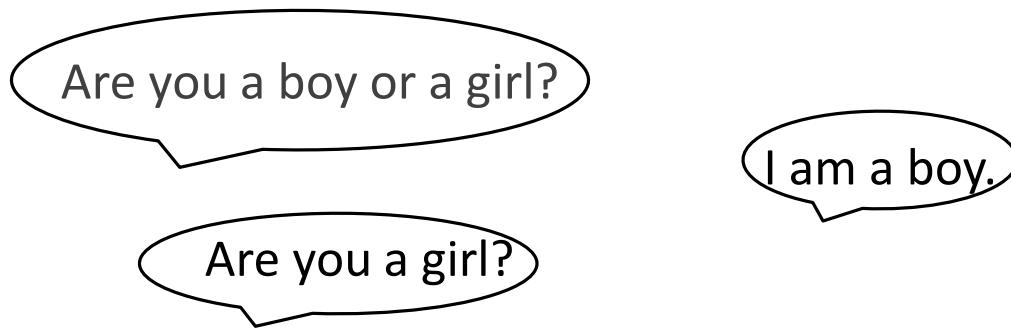
Are you a boy or a girl?

I am a boy.



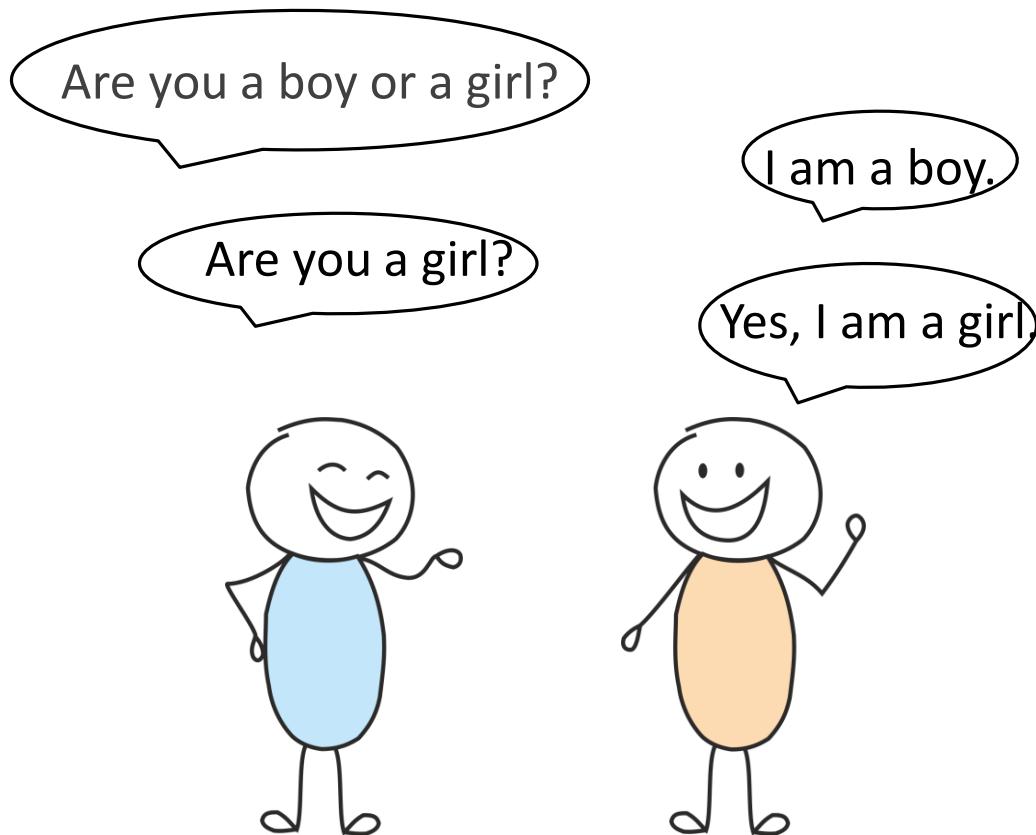
Storyline 3.

The importance of coherence.



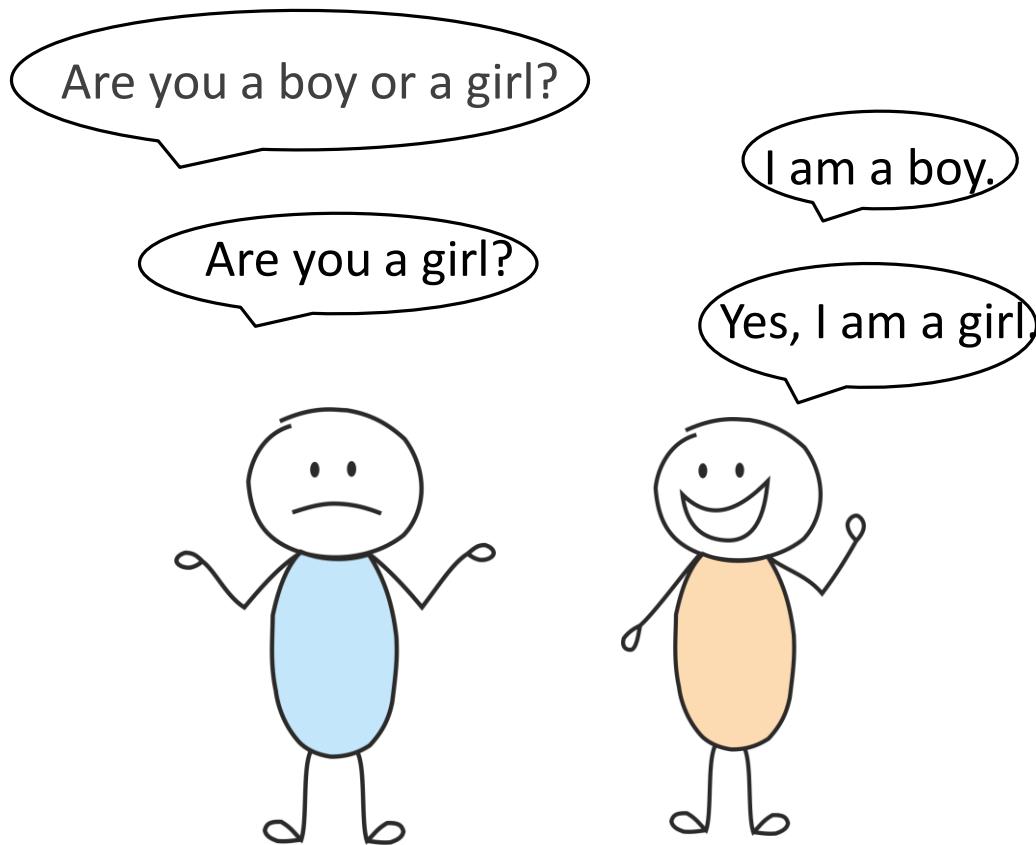
Storyline 3.

The importance of coherence.

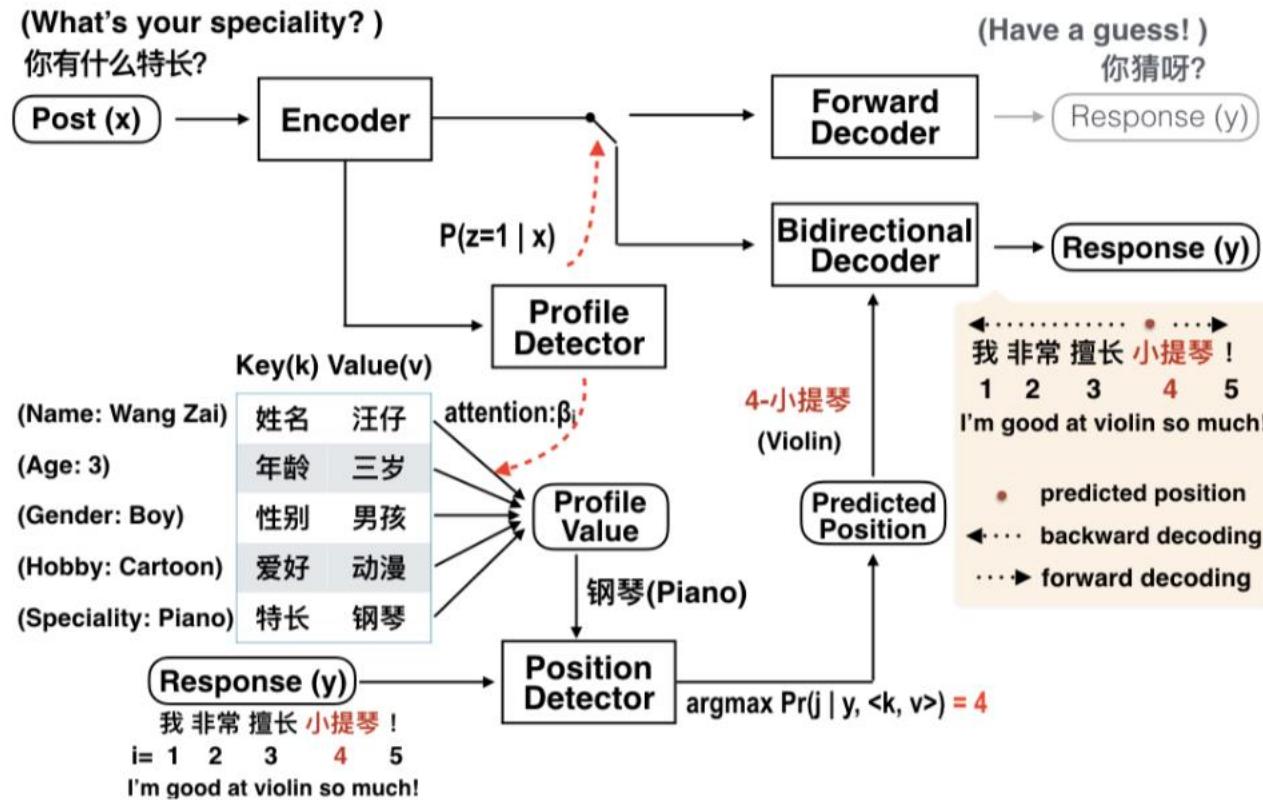


Storyline 3.

The importance of coherence.

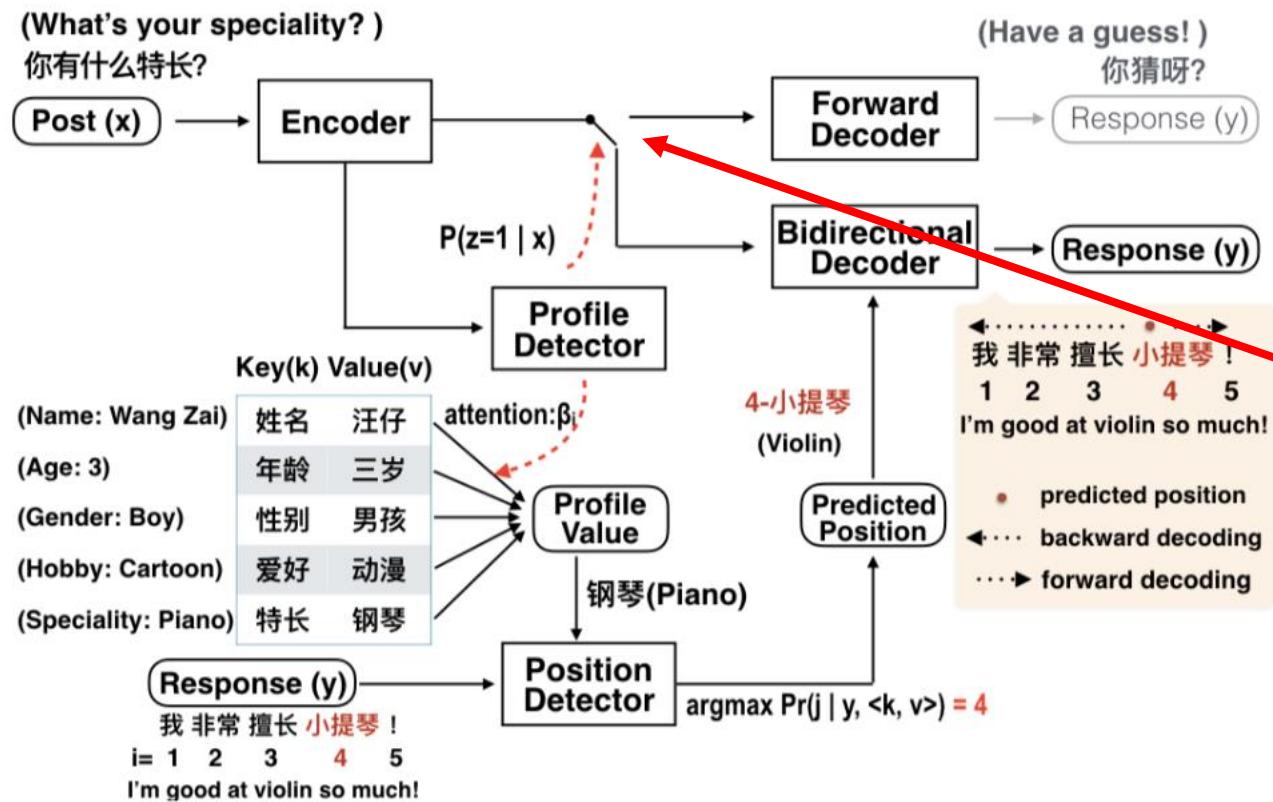


Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

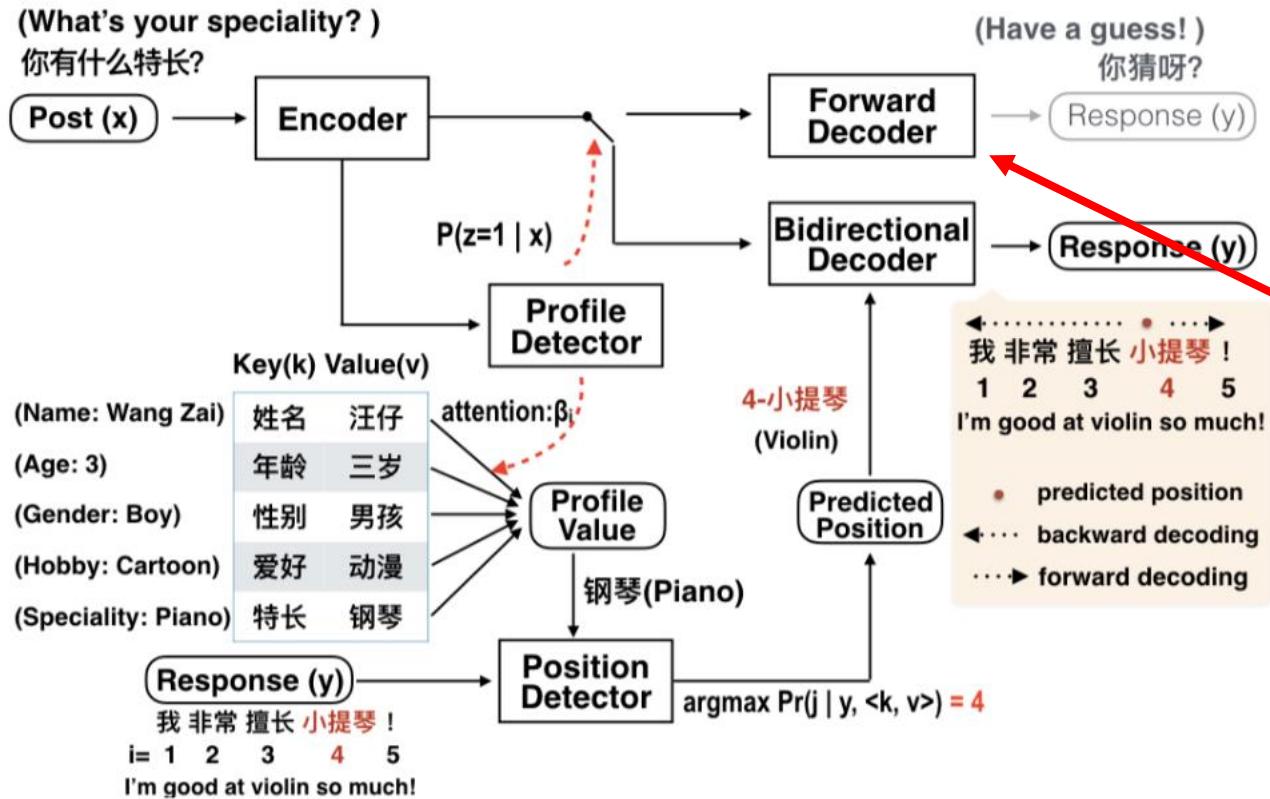
Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

Decide if we need profile

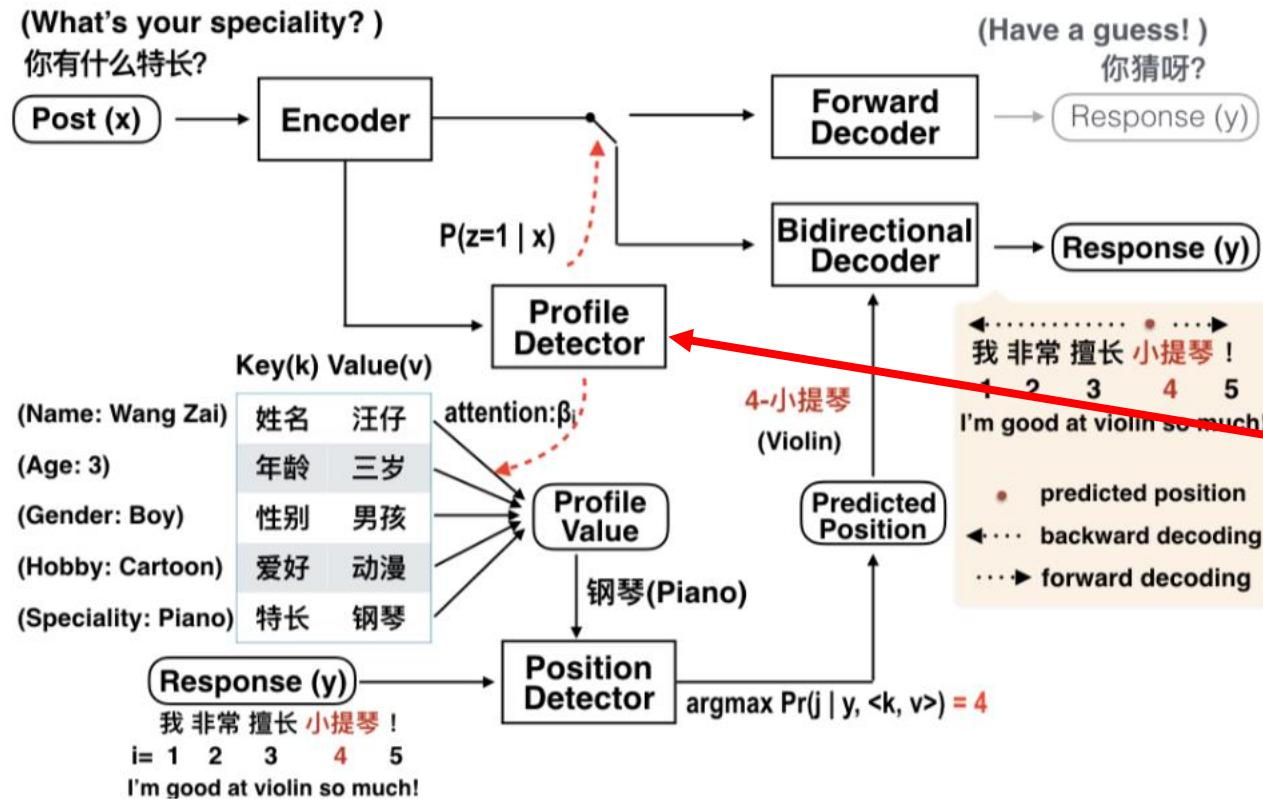
Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

If don't need,
use general decoder

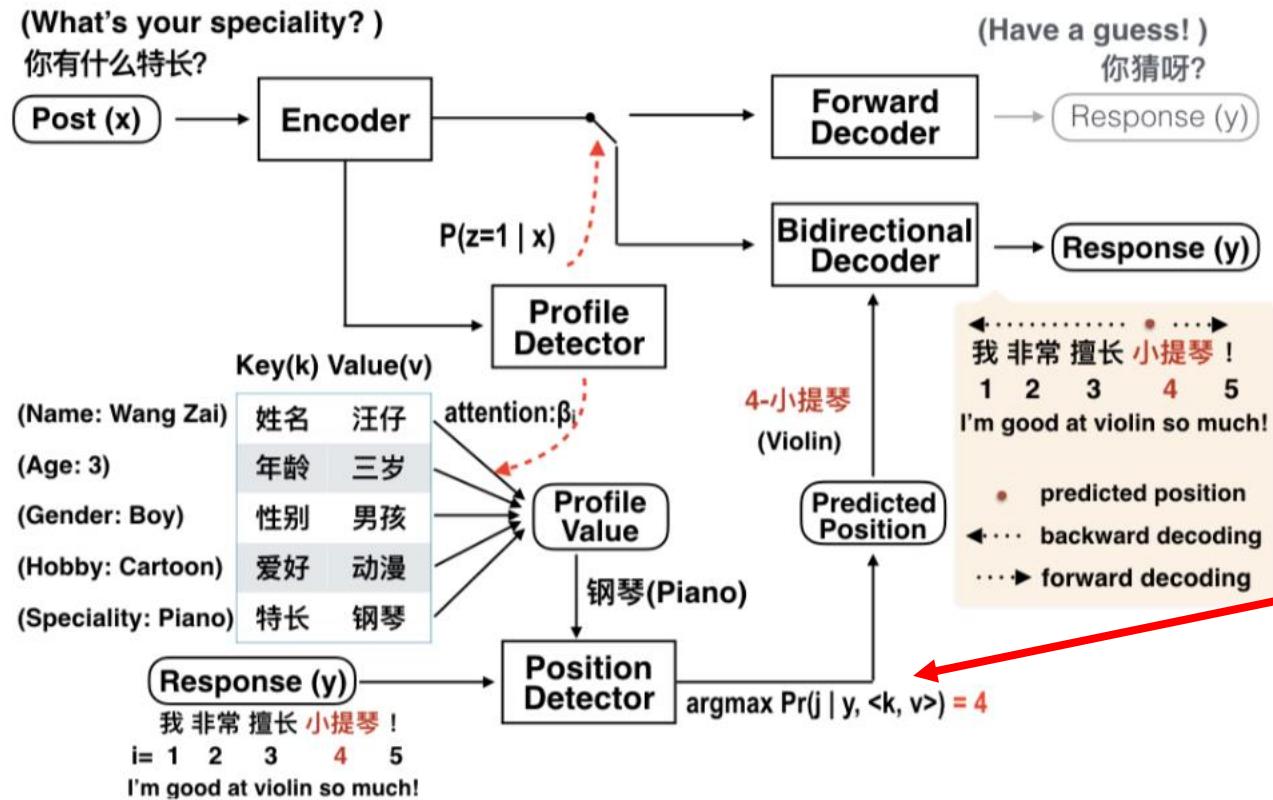
Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

If need, pick profile value

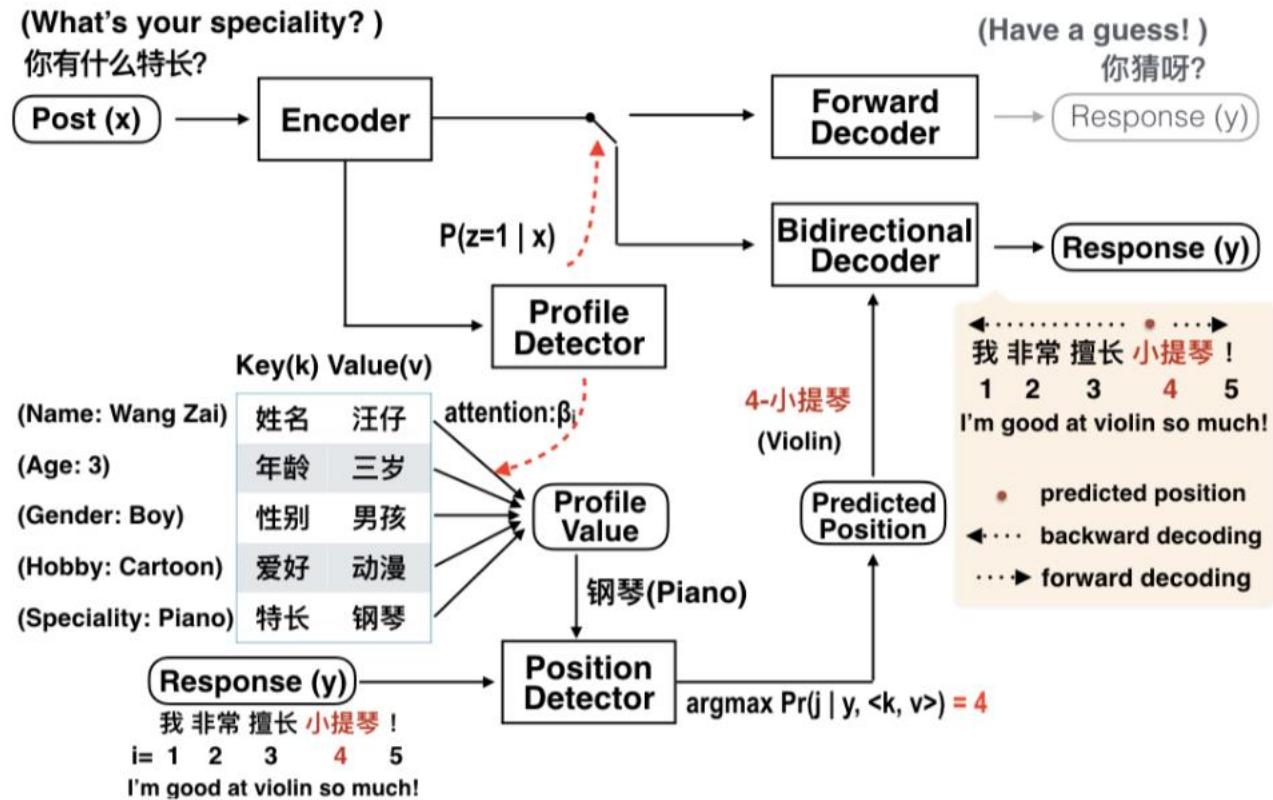
Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

Predict position for the
profile value

Assigning Personality/Identity to a Chatting Machine



Goal: to generate responses
that are coherent to its
prespecified
identity/personality

Generate from this position
both ways

Assigning Personality/Identity to a Chatting Machine

| Chinese | English(Translated) |
|----------------|---------------------------------|
| U:你还没说你几岁呢 | U:You haven't told me your age. |
| S:我三岁了 | S:I'm three years old. |
| U:你今年有15了不 | U:Are you 15 years old or not? |
| S:我还没到呢 | S:I'm not yet. |
| U:你多大啦 | U:How old are you? |
| S:3岁了 | S:Three years old. |

Table 5: Samples of consistent conversations generated by our model. *U/S* indicates User/System.

Goal: to generate responses that are coherent to its prespecified identity/personality

How NOT To Evaluate Your Dialogue System

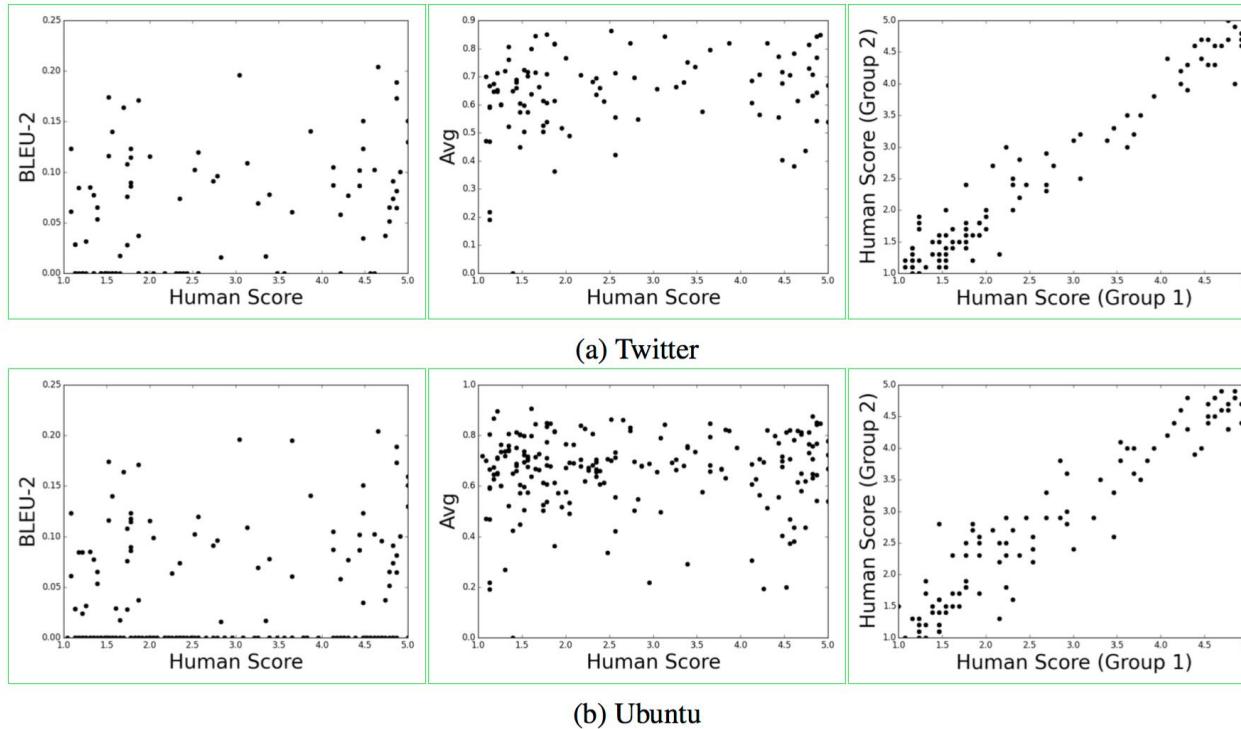


Figure 1: Scatter plots showing the correlation between metrics and human judgements on the Twitter corpus (a) and Ubuntu Dialogue Corpus (b). The plots represent BLEU-2 (left), embedding average (center), and correlation between two randomly selected halves of human respondents (right).

All metrics show either weak or no correlation with human judgements

How NOT To Evaluate Your Dialogue System

Context of Conversation

Speaker A: Hey John, what do you want to do tonight?

Speaker B: Why don't we go see a movie?

Ground-Truth Response

Nah, I hate that stuff, let's do something active.

Model Response

Oh sure! Heard the film about Turing is out!

Table 1: Example showing the intrinsic diversity of valid responses in a dialogue. The (reasonable) model response would receive a BLEU score of 0.

What else tried to improve?

- Improve diversity using latent variables
- Topic consistency
- Use outside knowledge base
- Learn through interaction
- Evaluation (BLEU, word2vec metrics are bad)

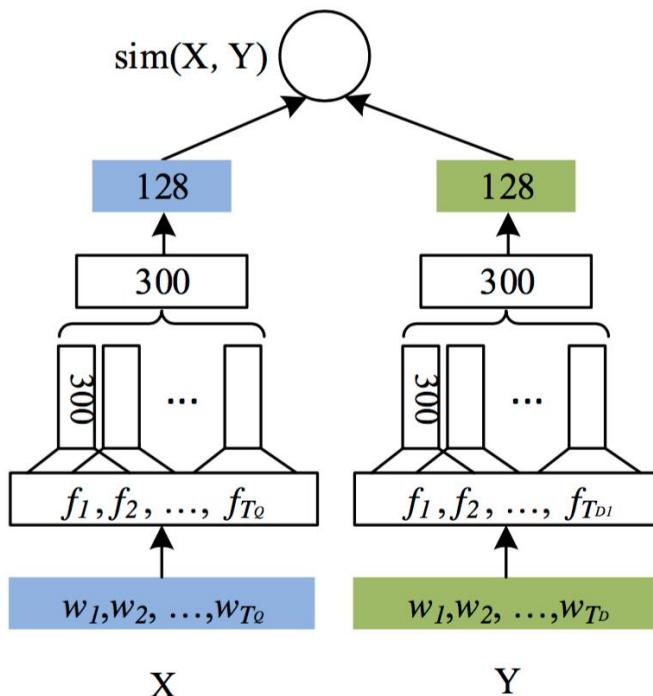
More details and links here: <https://arxiv.org/pdf/1711.01731.pdf>

Retrieval-based methods for General conversation



Deep Semantic Similarity Model (DSSM)

| | |
|-----------------------------------------|-------|
| Relevance measured by cosine similarity | |
| Semantic layer | h |
| Max pooling layer | ν |
| Convolutional layer | c_t |
| Word hashing layer | f_t |
| Word sequence | x_t |



Learning: maximize the similarity between X (source) and Y (target)

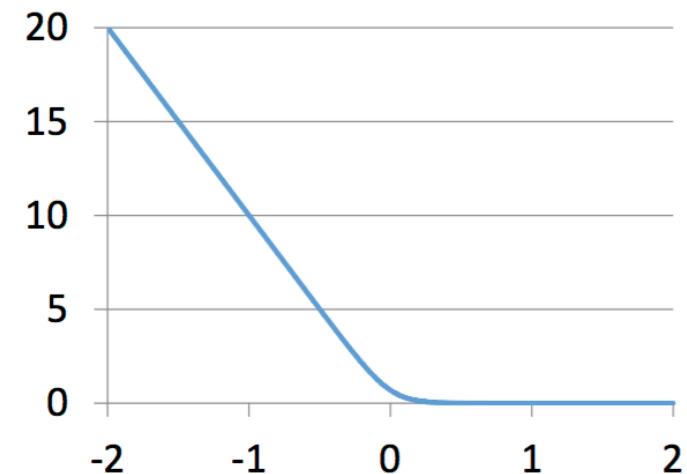
Representation: use DNN to extract abstract semantic representations

Convolutional and Max-pooling layer: identify key words/concepts in X and Y

Word hashing: use sub-word unit (e.g., letter n -gram) as raw input to handle very large vocabulary

Deep Semantic Similarity Model (DSSM)

- Consider a query X and two examples Y^+ and Y^-
Assume Y^+ is better than Y^- for query X
- $sim_{\theta}(X, Y^-)$ is the cosine similarity of X and Y in semantic space, mapped by DSSM parameterized by θ
- $\Delta = sim_{\theta}(X, Y^+) - sim_{\theta}(X, Y^-)$
- $Loss(\Delta, \theta) = \log(1 + \exp(-\gamma\Delta))$



What can be different?

- Ways to encode sentences (utterances, documents, etc.) – CNNs, RNNs, MLP, etc
- Loss functions
- Choice of number of context sentences and ways to encode them
- Negative sampling schedule (Yandex.Search experience)

But it's just a matching task!

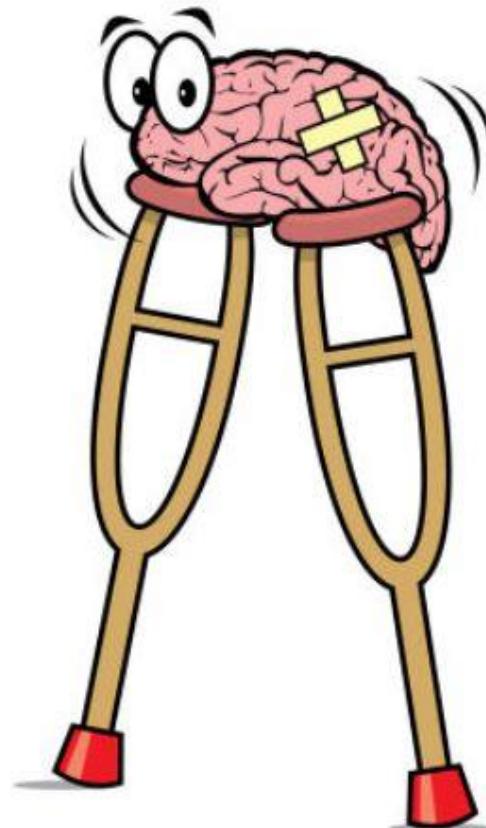
What can be different?

- Ways to encode sentences (utterances, documents, etc.) – CNNs, RNNs, MLP, etc
- Loss functions
- Choice of number of context sentences and ways to encode them
- Negative sampling schedule (Yandex.Search experience)

But it's just a matching task!

- web search query and clicked document
- utterance and response in dialogue
- user profile and article
- document and key phrases to be highlighted
- key phrase and entity with its wiki page

Hack of the day



Negative examples

Recall how we train DSSM:

- $\Delta = sim_{\theta}(X, Y^+) - sim_{\theta}(X, Y^-)$
- $Loss(\Delta, \theta) = \log(1 + \exp(-\gamma\Delta))$

We want a true answer to
be closer to query, than a
wrong one

How do we get Y^- ?

Negative examples

Recall how we train DSSM:

- $\Delta = sim_{\theta}(X, Y^+) - sim_{\theta}(X, Y^-)$
- $Loss(\Delta, \theta) = \log(1 + \exp(-\gamma\Delta))$

We want a true answer to
be closer to query, than a
wrong one

How do we get Y^- ?

Usually, it's a random example (remember for example word2vec?)

Problem

How do we get Y^- ?

Usually, it's a random example (remember for example word2vec?)

What if you have plenty of examples, and only a few good ones?

We want a true answer to be closer to query, than a wrong one



Problem

How do we get Y^- ?

Usually, it's a random example (remember for example word2vec?)

We want a true answer to be closer to query, than a wrong one

What if you have plenty of examples, and only a few good ones?

Would it be so hard to distinguish between a true answer and a random example?



Hard-Negatives!

Give as “negative” those negative examples, for
which your network gives high score
(thinks that they’re good)

Negative mining

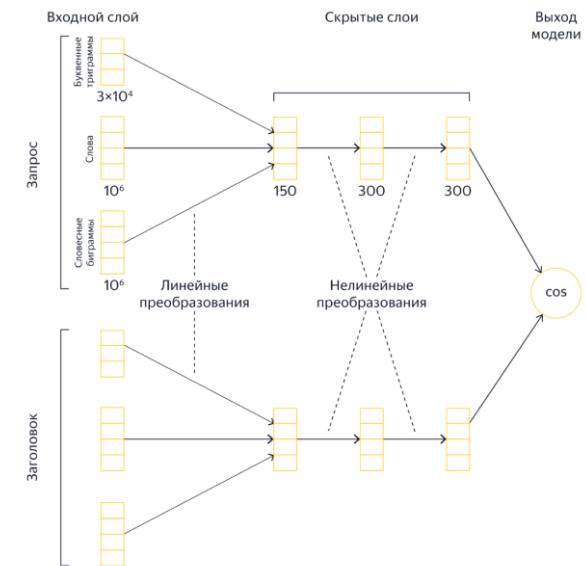
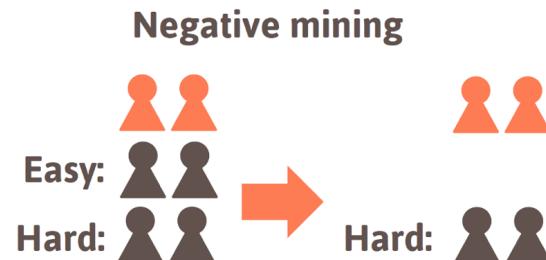


Hard-Negatives!

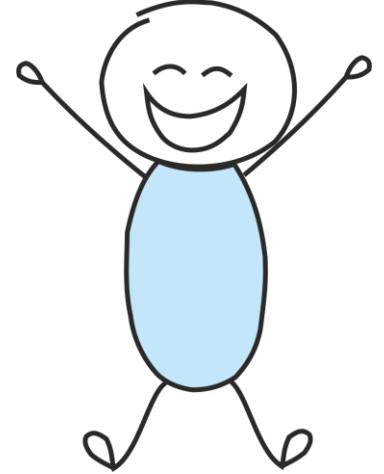
Give as “negative” those negative examples, for which your network gives high score (thinks that they’re good)

Real-life story: Yandex.Search and DSSM

<https://habr.com/company/yandex/blog/314222/>

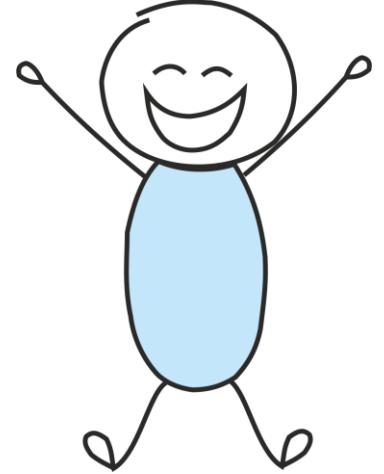


That's all for today!



Always yours,
Yandex Research

That's all for today!

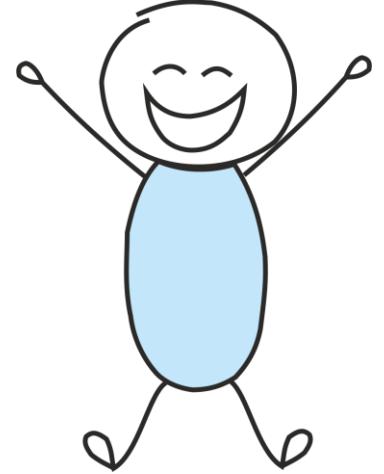


In the next episode:

- Adversarial methods in NLP

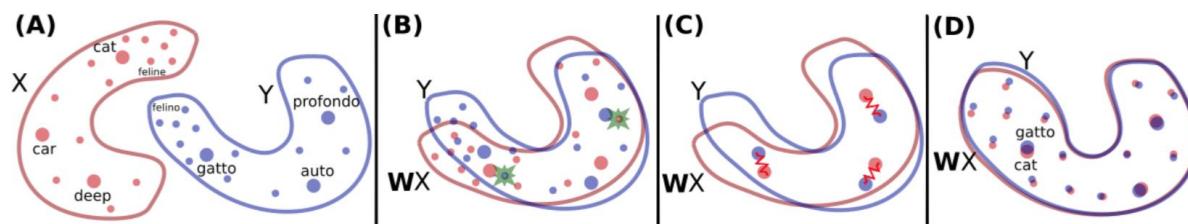
Always yours,
Yandex Research

That's all for today!



In the next episode:

- Adversarial methods in NLP
+ unsupervised embedding space mapping!



Always yours,
Yandex Research