# Conversational AI in Fiction

# Evolution of Conversational AI

Launch of Apple Siri, 2011

World Wide Web - largest knowledge base, 1989

The Dartmouth proposal, 1955

AT&T MATCH Multimodal System, 2002

Future of CAIs
- Integration to many services
- Multimodal
- Multilingual
- Multi-turn
- Ubiquitous

The ELIZA chatbot, 1964

Virtual Assistant platforms, 2017

Turing Imitation Game, 1950
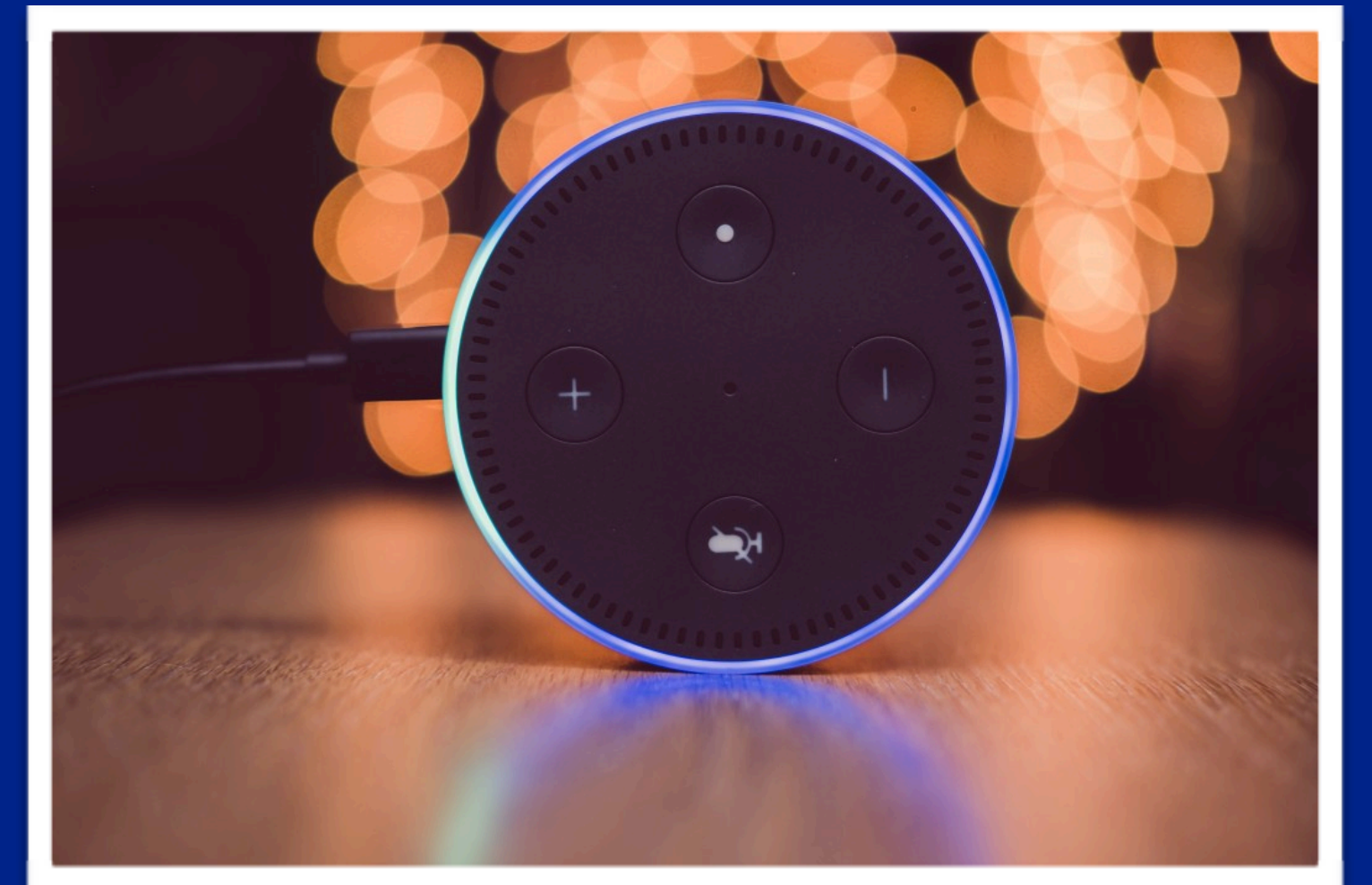
POLYAi

# A Not-So-Far-Away Future

*"Hey PolyBot, could you order a pizza home for me. I will be back around 7.30. You know my preference, right?"*
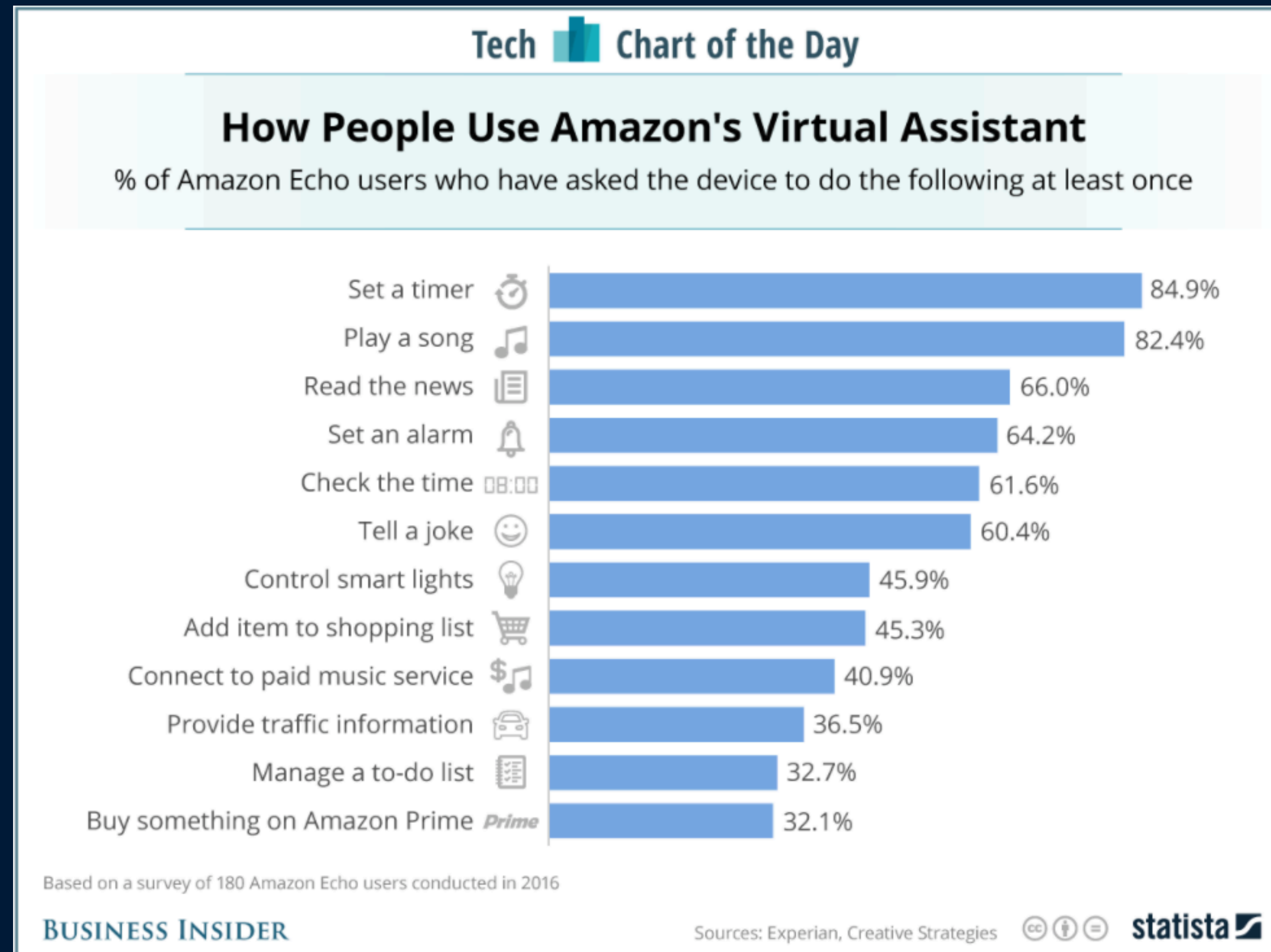
*"Hey PolyBot, what is the status of my order? Could you tell PizzaHut that I'm going to be 30 minutes late?"*

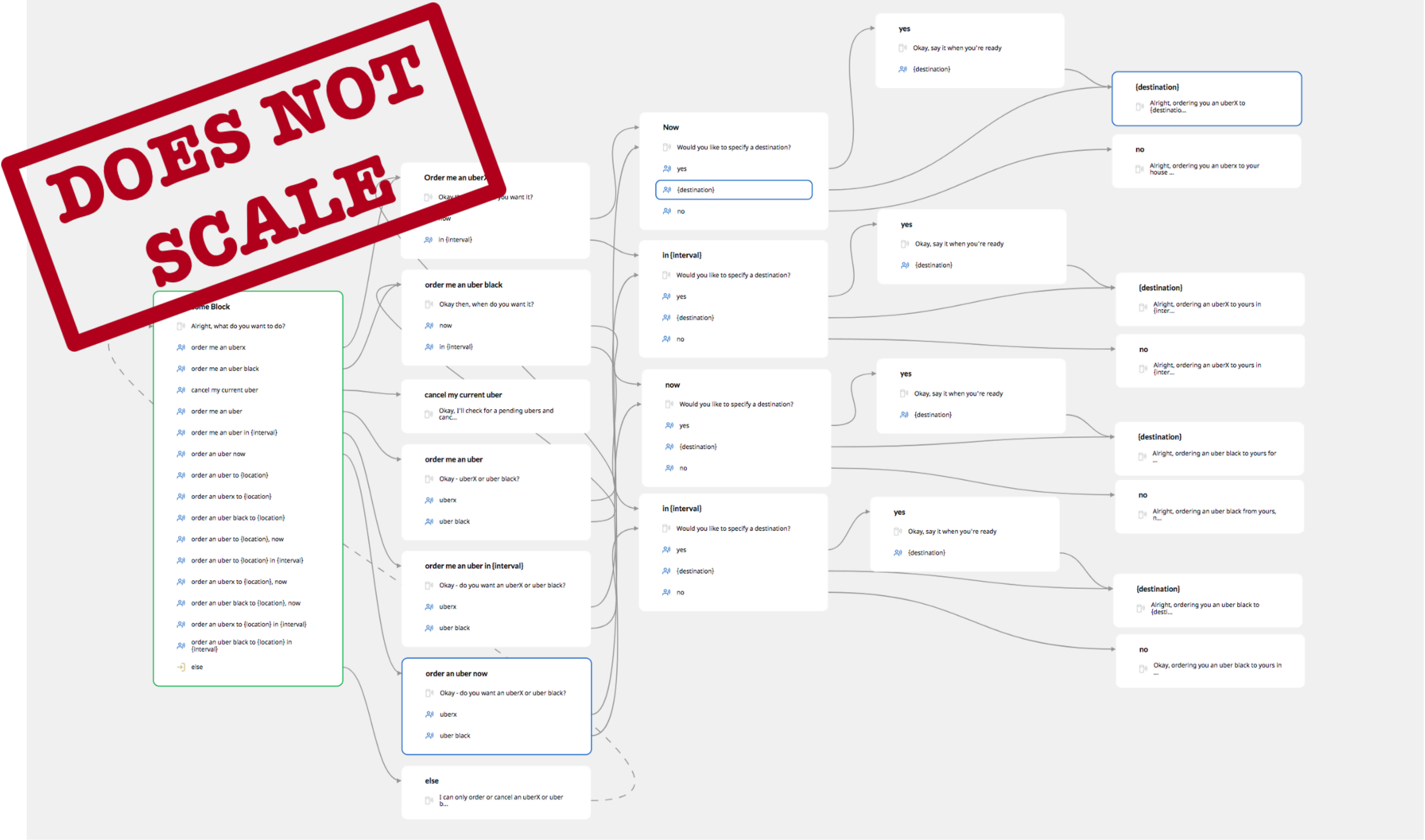*"[Notification] Your Margarita order has arrived. Please give a 1-5 star feedback on our delivery service."*

POLYAI

# Conversational AI in 2018



Tech 📊 Chart of the Day

**How People Use Amazon's Virtual Assistant**

% of Amazon Echo users who have asked the device to do the following at least once

| Action | % |
|---|---|
| Set a timer | 84.9% |
| Play a song | 82.4% |
| Read the news | 66.0% |
| Set an alarm | 64.2% |
| Check the time | 61.6% |
| Tell a joke | 60.4% |
| Control smart lights | 45.9% |
| Add item to shopping list | 45.3% |
| Connect to paid music service | 40.9% |
| Provide traffic information | 36.5% |
| Manage a to-do list | 32.7% |
| Buy something on Amazon Prime | 32.1% |

Based on a survey of 180 Amazon Echo users conducted in 2016

**BUSINESS INSIDER**

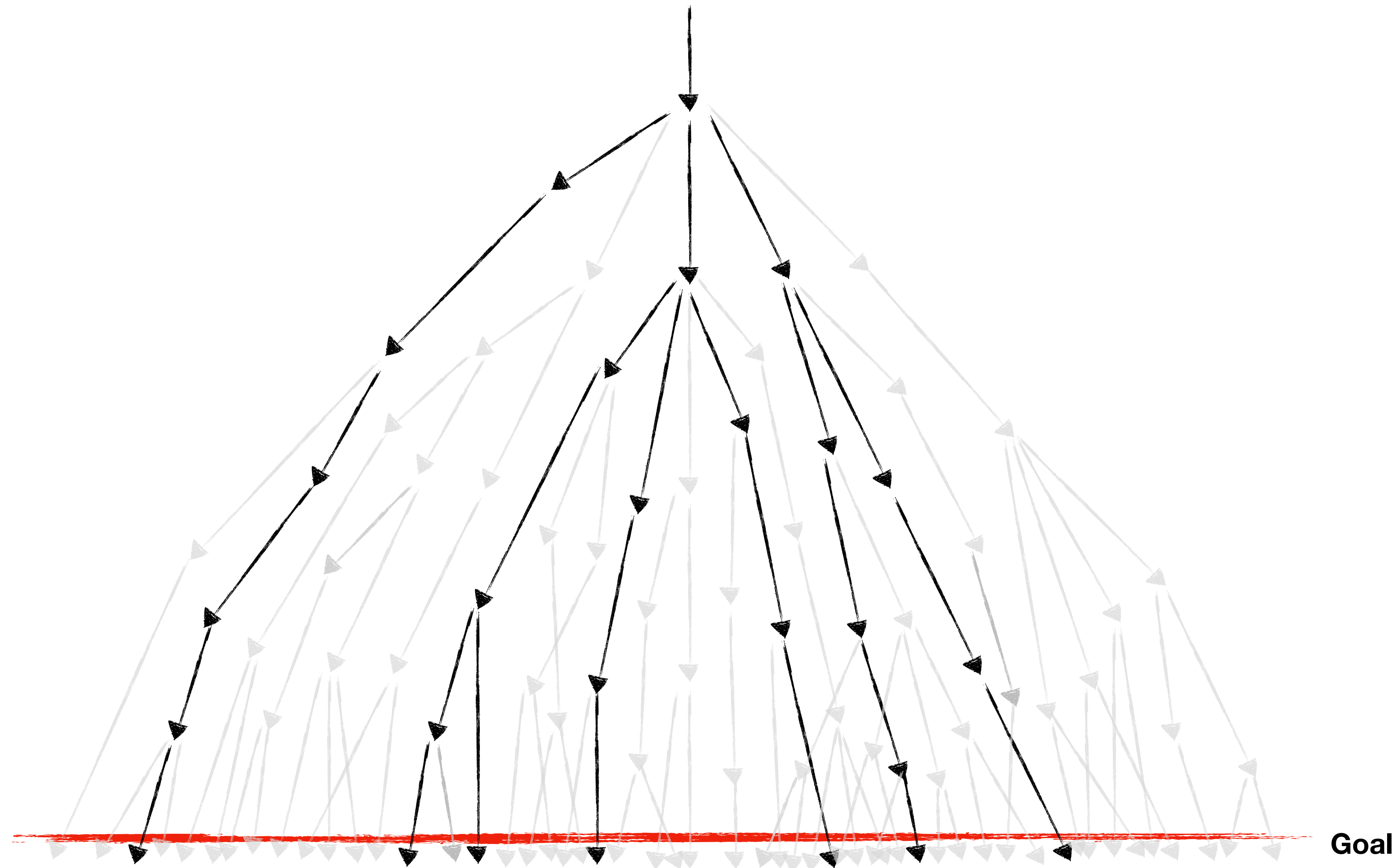Sources: Experian, Creative Strategies

statista

POLYAI

# Conversational AI in 2018

- Virtual Assistants are currently used for **simple tasks** - the more complex (and useful) ones are implemented and supported by the platform providers themselves

- Third-party service providers are finding it **hard to connect** their services to existing virtual assistant platforms

  - This major **technical block** stands in the way of virtual assistants becoming the next iPhone - they are waiting for their AppStore to come to life

POLYAi

# Current Industry Practice - Flowcharts

POLYAi

# The Happy Path Approach



Goal

POLYAi

- Set an alarm for 8.30am.

- Alarm at 8.30am, please.

- Wake me up in seven hours.

- Alarm next morning at 8.30.

- Wake me up at 8.30am.

- An alarm for half eight.

- Turn on my 8.30 alarm.

- …

POLYAI

# Simple Flowchart Tools = Simple Voice-Based Apps

## Top Alexa Skills are still simple single-turn commands



"Alexa, ask Question of the Day."

**Question of the Day**
⭐⭐⭐⭐⭐ 684

"Alexa, launch National Rail."

**National Rail**
⭐⭐☆☆☆ 420

"Alexa, play 'Would you rather?'"

**Would You Rather?**
⭐⭐⭐⭐½ 271

"Alexa, start five minute plank"

**5-Minute Plank Workout**
⭐⭐⭐⭐½ 43

"Alexa, ask Jurassic Bark to bark"

**Jurassic Bark**
⭐⭐⭐⭐☆ 13

"Alexa, ask Recipedia for a quick chicken recipe."

**Recipedia**
⭐⭐⭐☆☆ 63

"Alexa, open the BBC"

**BBC**
⭐⭐⭐☆☆ 216

"Alexa, play river sounds"

**Ambient Noise: River Sounds**
⭐⭐⭐⭐☆ 6

"Alexa, open KISS FM"

**KISS FM**
⭐⭐⭐☆☆ 11

"Alexa, open Smooth Radio"

**Smooth Radio**
⭐⭐⭐⭐⭐ 11

"Alexa, meow meow"

**Meow!**
⭐⭐⭐⭐☆ 95

"Alexa, what's my Flash Briefing?"

**BBC World Service**
⭐⭐⭐½☆ 30

"Alexa, Ask For A Fart"

**4AFart**
⭐⭐⭐⭐½ 61

"Alexa, what's my Flash Briefing?"

**BBC World Service**
⭐⭐⭐½☆ 30

"Alexa, play river sounds"

**Ambient Noise: River Sounds**
⭐⭐⭐⭐☆ 6

POLYAi

# Conversational AI Approaches

Unstructured

Pattern Matched-
based Chatbots

Retrieval Based
Chatbots

Neural Generative
Chatbots

Rule-based
Dialogue Systems

POMDP-based
Dialogue Systems

Neural Generative
Dialogue Systems

Structured

# of Data Available

POLYAi

# A Single-Direction Development?

Expert System → Bayesian Methods → Deep Learning → Latent Variable Models → ?

*"Is new technology always better than what we have already had?"*

POLYAi

# A Collective Effort!



Realizing Conversational AI requires efforts from multiple parties:

- 3rd-Party Service Providers.
- Conversational UX/UI Developers.
- Machine Learning Scientists/Engineers.

Realizing Conversational AI also requires a hybrid of approaches to tackle different parts of the pipeline:

- Expert programmed algorithms to handle domain-specific **business logics**.
- Different machine learning models to handle **prediction** and **uncertainty estimation**.
- Integration to **3rd-party services** to provide meaningful applications.
- Access to data to learn from human interactions.

POLY**Ai**

# Outline

- Conversational AI Architecture

- Overview of Conversational AI Research

  - End-to-End Chatbots

  - Learning-based Goal-oriented Systems

- Data! Data! Data!

- Conclusion

POLYAi

# Conversational AI - Architecture



"Book me a cab to Russell Square"

Speech Recognition

speech

Language Understanding

text

semantics

inform(service=taxi, dest=Russell Square)

Dialogue Management

Third Party APIs

speech

text

semantics

request(depart_time)

Speech Synthesis

Response Generation

"When do you want to leave?"

16

POLYAi

# Conversational AI - Architecture



"Book me a cab to Russell Square"

speech → Speech Recognition → text → Language Understanding → semantics → inform(service=taxi, dest=Russell Square)

Dialogue Management ↔ Third Party APIs

Response Generation → semantics → request(depart_time)

speech ← Speech Synthesis ← text ← Response Generation

"When do you want to leave?"

POLYАi

# Conversational AI - Architecture

# Conversational AI - Architecture



"Book me a cab to Russell Square"

Speech Recognition → Language Understanding → inform(service=taxi, dest=Russell Square) → Dialogue Management → Third Party APIs

speech / text / semantics

Dialogue Management → request(depart_time) → Response Generation → Speech Synthesis

semantics / text / speech

"When do you want to leave?"

POLYAi

# Conversational AI - Architecture



*"Book me a cab to Russell Square"*

Speech Recognition → Language Understanding

*speech* → *text* → *semantics*

inform(service=taxi, dest=Russell Square)

Dialogue Management ↔ Third Party APIs

*speech* ← *text* ← *semantics*

Speech Synthesis ← Response Generation

request(depart_time)

*"When do you want to leave?"*

POLYAi

# Sequence-to-Sequence Conversation Model

*Vinyals & Le, 2015*



**Human:** *what do you think about tesla ?*
**Machine:** *he 's a good conductor .*
**Human:** *what do you think about bill gates ?*
**Machine:** *he 's a good man .*
**Human:** *What do you think about messi ?*
**Machine:** *he 's a great player .*
**Human:** *what do you think about cleopatra ?*
**Machine:** *oh , she 's very regal .*
**Human:** *what do you think about england during the reign of elizabeth ?*
**Machine:** *it was a great place .*

**Human:** *what is the purpose of life ?*
**Machine:** *to serve the greater good .*
**Human:** *what is the purpose of living ?*
**Machine:** *to live forever .*
**Human:** *what is the purpose of existence ?*
**Machine:** *to find out what happens when we get to the planet earth .*
**Human:** *where are you now ?*
**Machine:** *i 'm in the middle of nowhere .*
**Human:** *what is the purpose of dying ?*
**Machine:** *to have a life .*
**Human:** *what is the purpose of being intelligent ?*
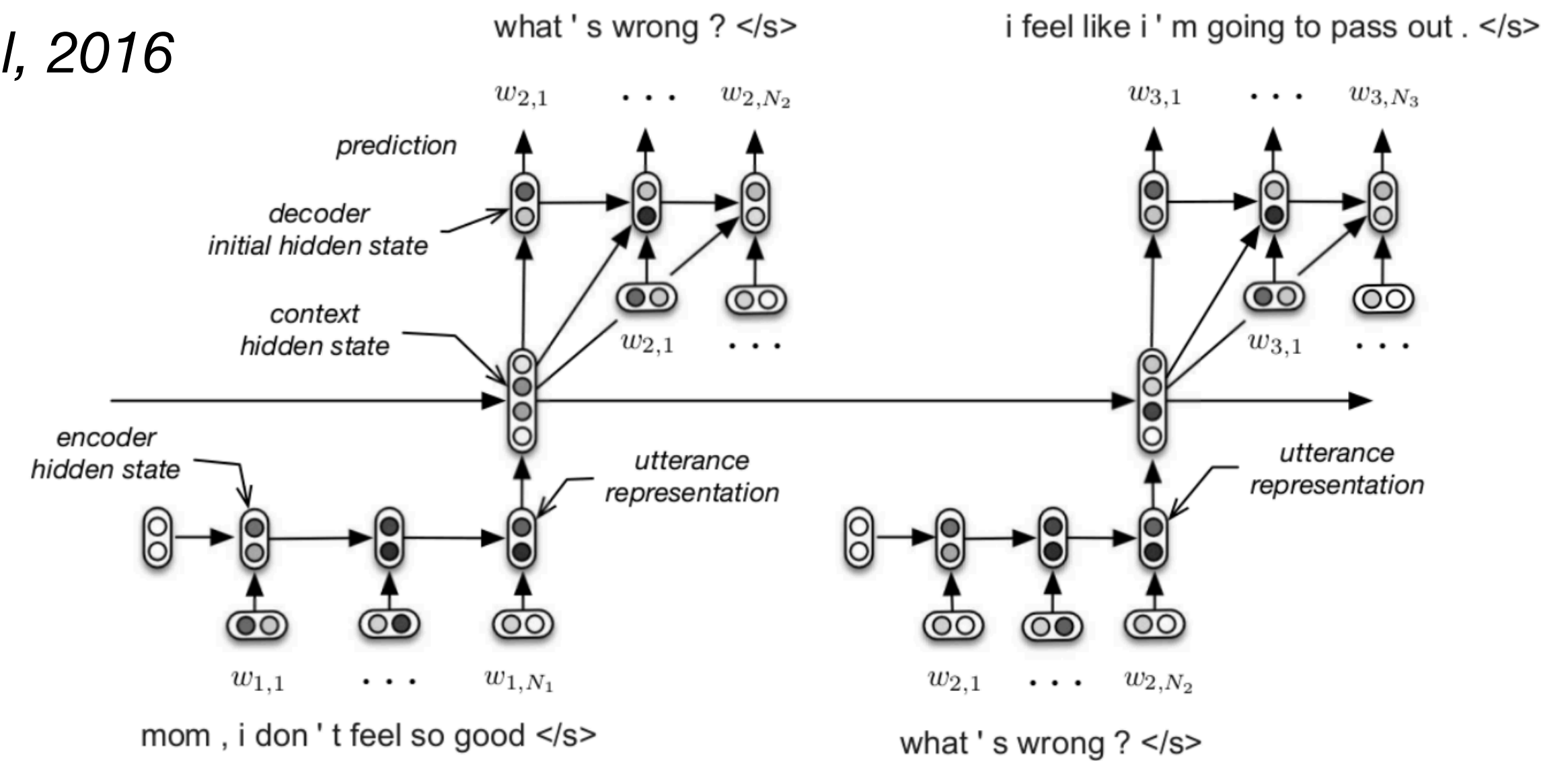**Machine:** *to find out what it is .*
**Human:** *what is the purpose of emotions ?*
**Machine:** *i don 't know .*

- Ignore the structure of the problem.

- A simple sequence-to-sequence model to learn a context-to-response mapping.

- Maximising a conditional probability p(response|context)

- Large model with huge data.

- Can learn to generate interesting responses.

POLYAi

# Hierarchical RNN for Dialogue
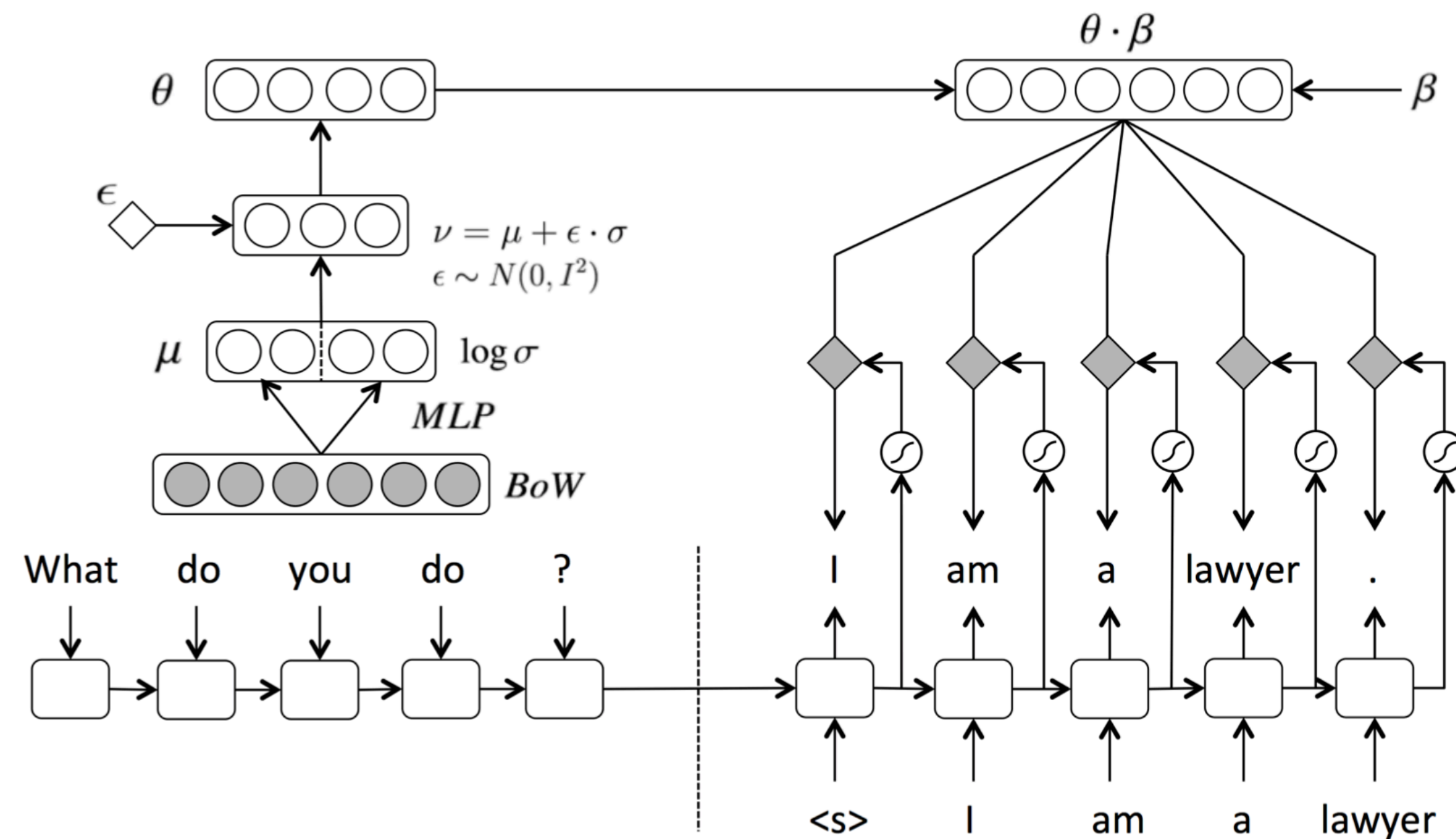
*Serban et al, 2016*



| **Reference** ($\mathbf{U_1}$, $\mathbf{U_2}$) | **MAP** | **Target** ($\mathbf{U_3}$) |
|---|---|---|
| $U_1$: yeah , okay . <br> $U_2$: well , i guess i ' ll be going now . | i ' ll see you tomorrow . | yeah . |
| $U_1$: oh . \<continued_utterance\> oh . <br> $U_2$: what ' s the matter , honey ? | i don ' t know . | oh . |
| $U_1$: it ' s the cheapest . <br> $U_2$: then it ' s the worst kind ? | no , it ' s not . | they ' re all good , sir . |
| $U_1$: \<person\> ! what are you doing ? <br> $U_2$: shut up ! c ' mon . | what are you doing here ? | what are you that crazy ? |

POLYAi

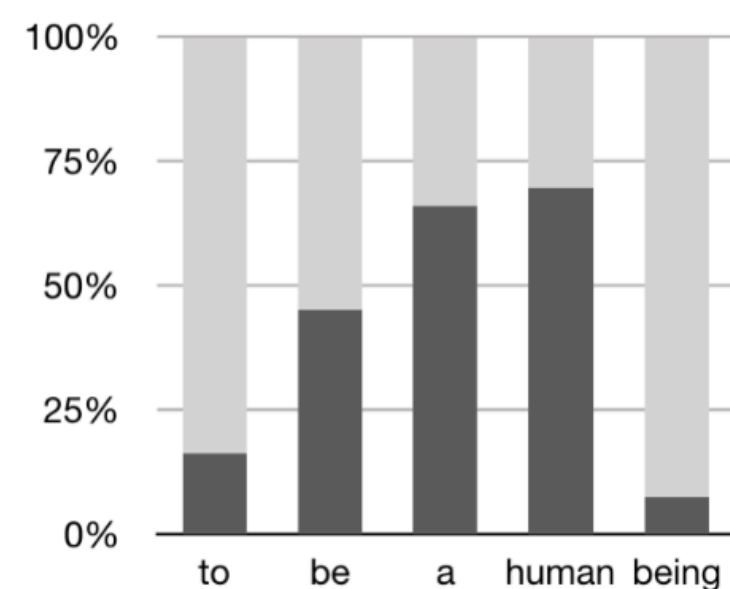# Problems of the Approach

- Problems

  - Generic, incoherent, and non-causal responses.

  - The goal is to "*response like a human being*", not very useful for business use cases.

  - Lack of integration to knowledge base or 3rd party services.

  - Evaluation is super subjective and difficult.

POLYAi

# An Improvement - Latent Topic Conversational Model



$\theta \cdot \beta$

$\theta$

$\epsilon$

$\nu = \mu + \epsilon \cdot \sigma$
$\epsilon \sim N(0, I^2)$

$\mu$     $\log \sigma$

$MLP$

$BoW$

$\beta$

What do you do ?

I am a lawyer .

<s> I am a lawyer

**Prompt**: What is the meaning of life?

**Prompt**: What do you think about Messi?

to be a human being

i think messi is the best player in the world
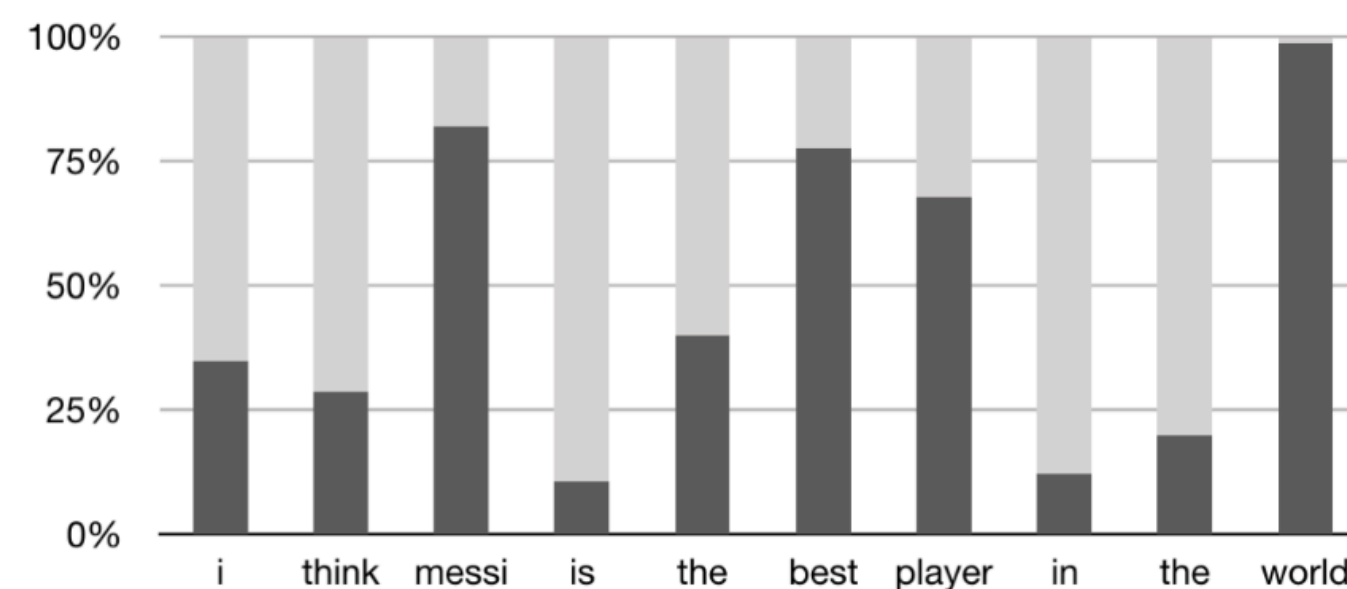
■ Topic contribution    ■ Seq2Seq contribution

- Combine topic and seq2seq models.

- Optimise as a Variational Autoencoder.

- Use the topic component to control meaning, while adopt seq2seq to produce grammatical responses.

**Prompt**: What is the purpose of existence?
**LTCM**: To be able to see the world.
To be able to see things.
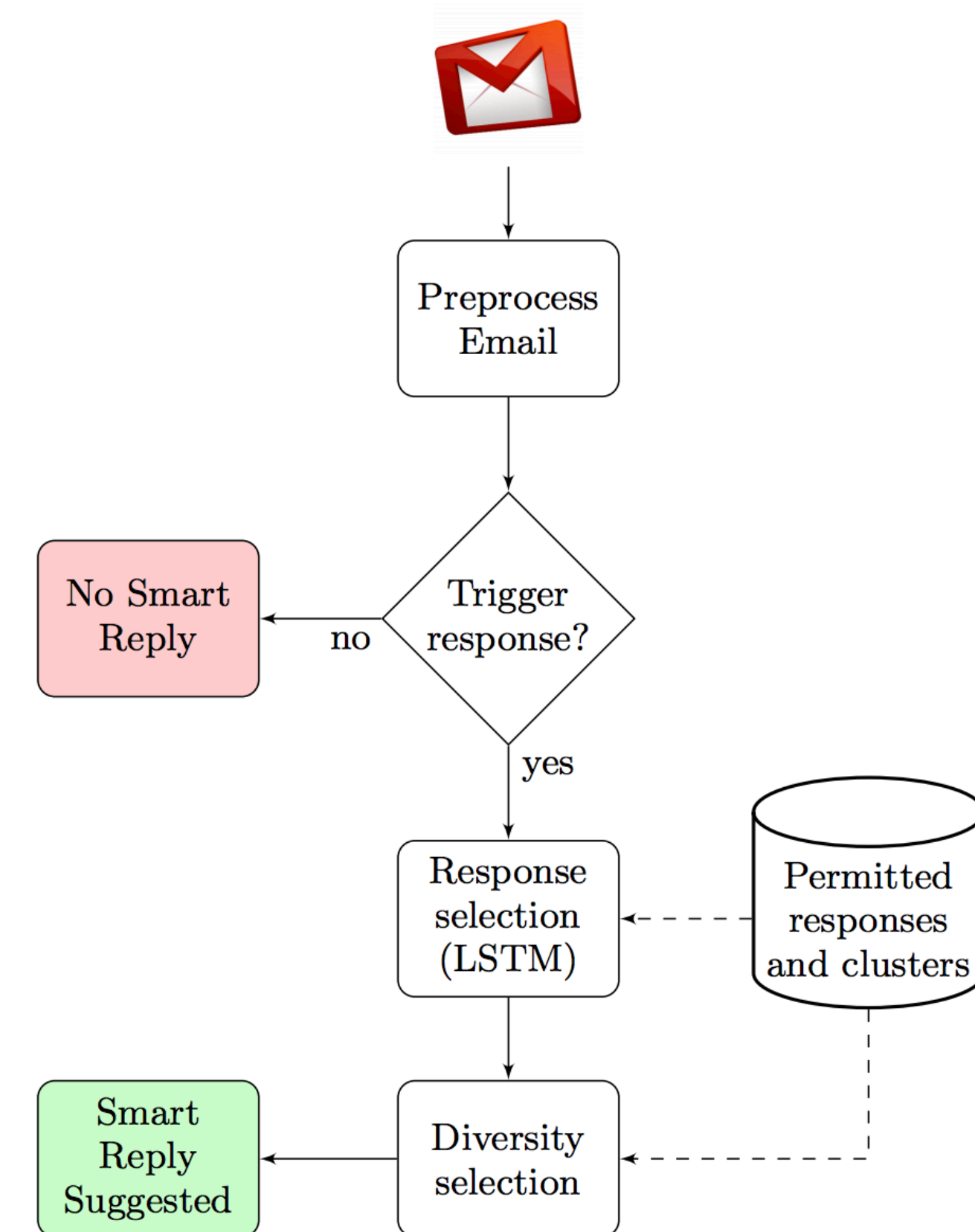To make it better.

*Wen et al, 2018 (to be appear)*

POLYAi

# Case Study: Google Smart Reply

- Response selected from a whitelist.

- LSTM is used for sentence scoring:

$$P(r_1, ..., r_m | o_1, ..., o_n) = \prod_{i=1}^{m} P(r_i | o_1, ..., o_n, r_1, ..., r_{i-1})$$
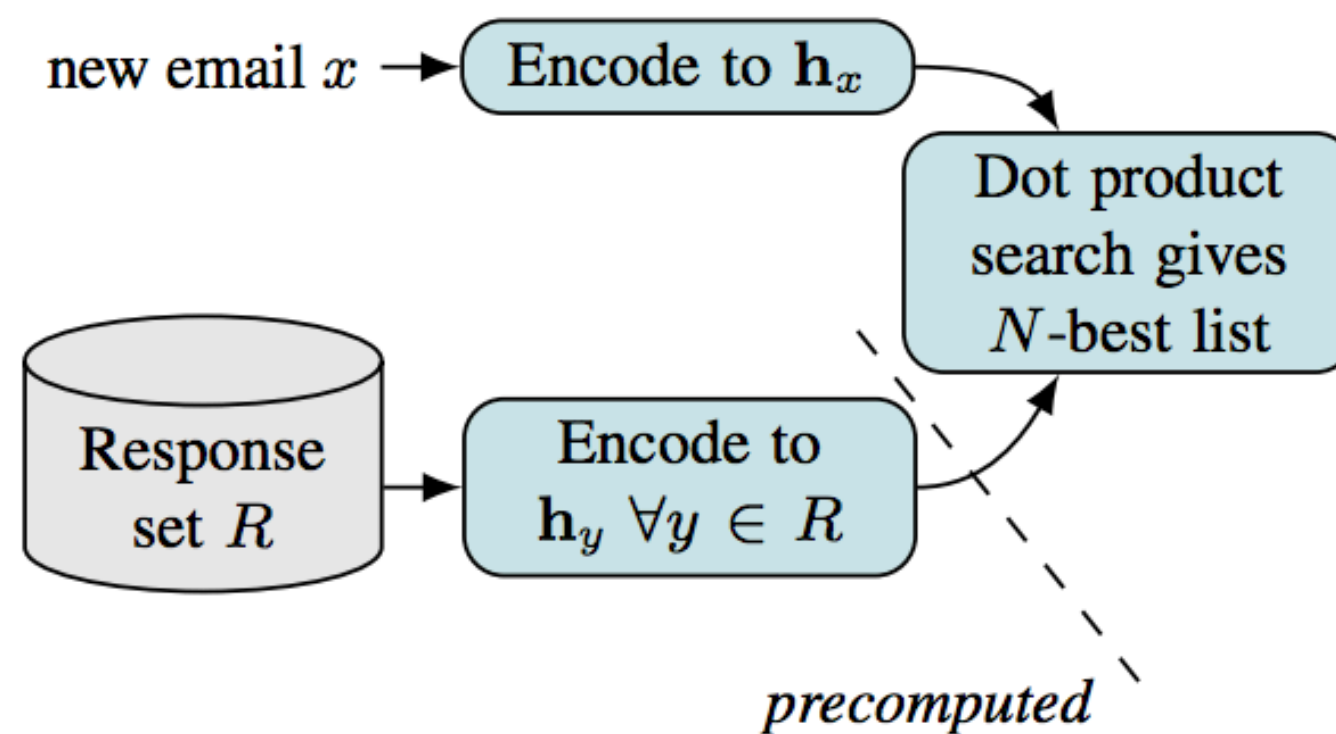
- Following the Seq2Seq paradigm.

- Inference is done by:

  - Organizing responses as a trie

  - Beam search

Kannan et al, 2016

# Case Study: Google Smart Reply

- Encode response set (precomputed)

- Dot product (fast) between email and response vectors.

- Hierarchical quantization for search.



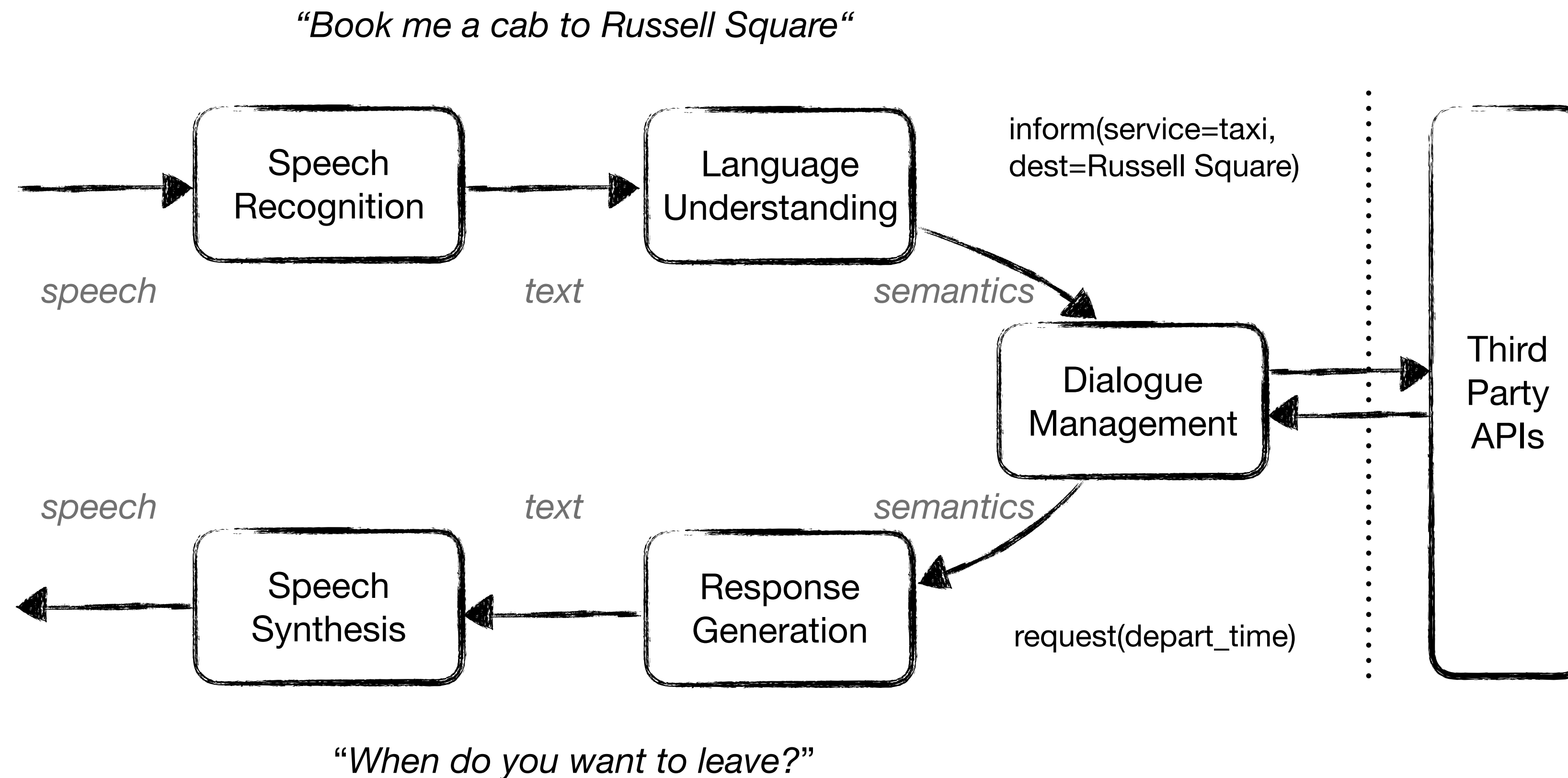| | Conversion Rate (relative to Seq2Seq) | Speed (relative to Seq2Seq) |
|---|---|---|
| Seq2Seq + beam search | 100% | 1x |
| Dot Product + hierarchical search | 104% | **100x** |

Henderson et al, 2017

# Short Summary

- Generative Sequence-to-Sequence model is extremely hard to get it to work.

  - It struggled in generating diverse responses - but other methods have been proposed to address this problem.

  - Generating the intended meaning while learning the semantic representation on the fly is still impossible.

- Runtime-wise, generative model is much slower than simple dot product retrieval-based model.

- Dot product model has been successfully deployed and shows a better result than generative models.
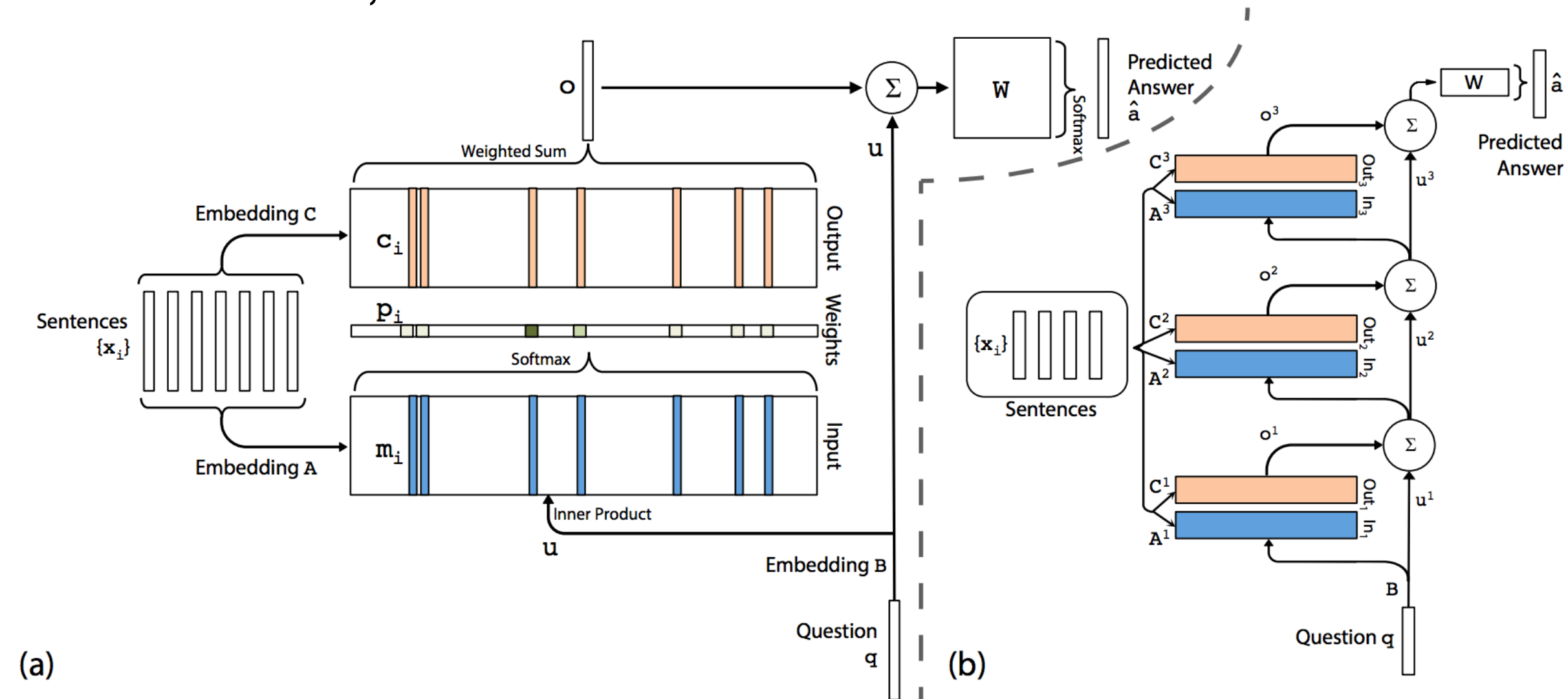
POLYAi

# Learning-based Goal-oriented Dialogue

# Recap - Conversational AI Architecture

"*Book me a cab to Russell Square*"

inform(service=taxi,
dest=Russell Square)

Speech
Recognition

Language
Understanding

*speech*

*text*

*semantics*

Dialogue
Management

Third
Party
APIs

*speech*

*text*

*semantics*

Speech
Synthesis

Response
Generation

request(depart_time)

"*When do you want to leave?*"

POLYAi

# Memory Network for Dialogue Modeling

Bordes et al, ICLR 2017



(a) (b)

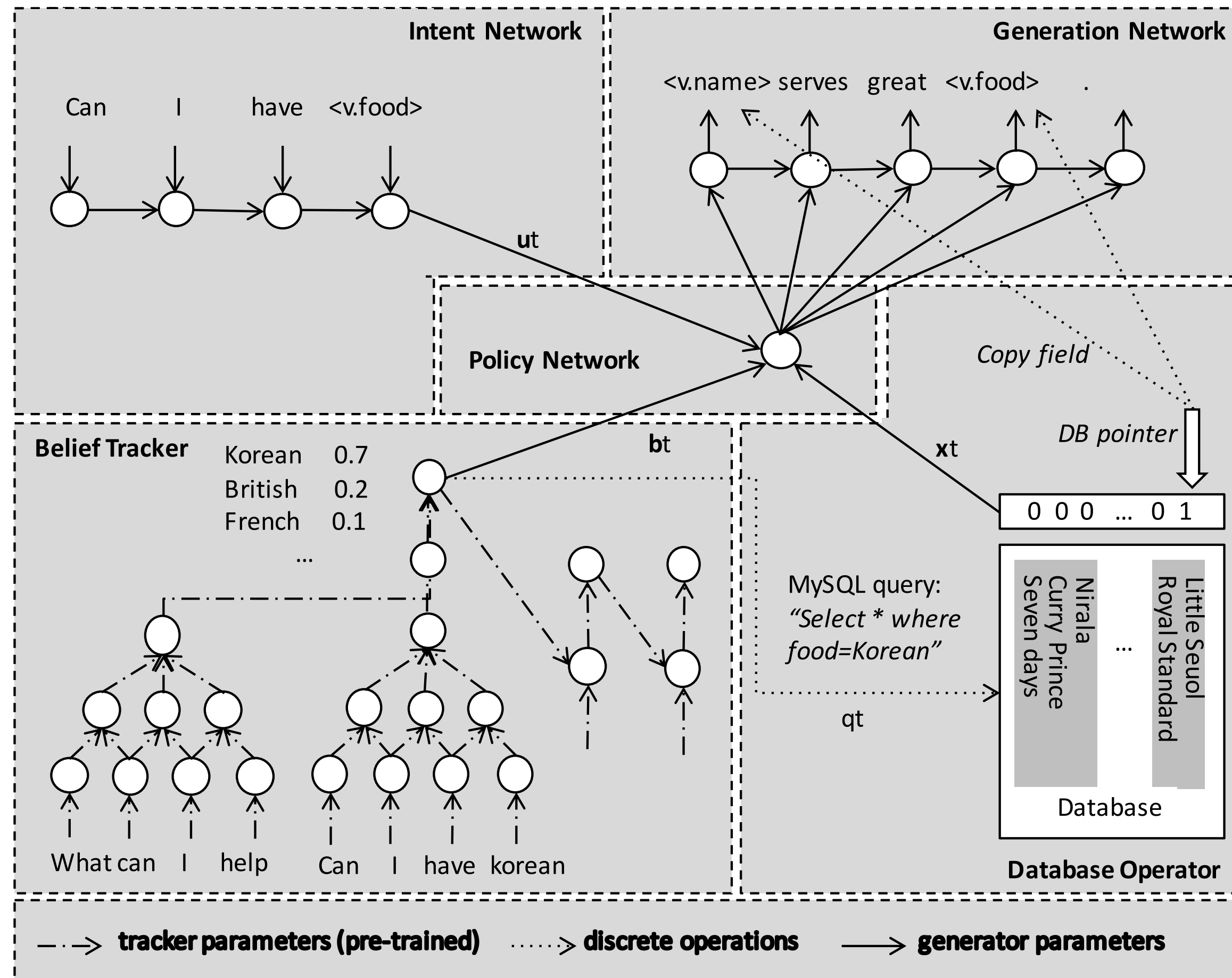| Time | Locutor | Dialog History | Hop #1 | Hop #2 |
|------|---------|----------------|--------|--------|
| 1 | User | hey concierge | .189 | .095 |
| 2 | User | could you check if i can get a rservation at \<org\> \<date\> for brunch | .209 | .178 |
| 3 | User | \<number\> people | .197 | .142 |
| 4 | User | \<silence\> | .187 | .167 |
| 5 | Bot | hi \<person\> unfortunately \<org\> is fully booked for \<date\> and there's \<number\> people on the waiting list | .225 | .410 |
| **User input** | | when's the earliest availability | | |
| **Correct answer** | | i'll check | | |
| **Pred. answer #1** | | i'm on it | [Incorrect] | |
| **Pred. answer #2** | | i'll find out | [Incorrect] | |
| **Pred. answer #3** | | i'll take a look | [Incorrect] | |
| **Pred. answer #4** | | i'll check | **[Correct]** | |
| **Pred. answer #5** | | i'll check into it | [Incorrect] | |

- Without using the structure of dialogue systems.

- Apply attention mechanism to search for relevant information in the context.

- Aggregate the current user query with the retrieved context to make predictions.

- Predictions include API calls and response selection from a candidate pool.

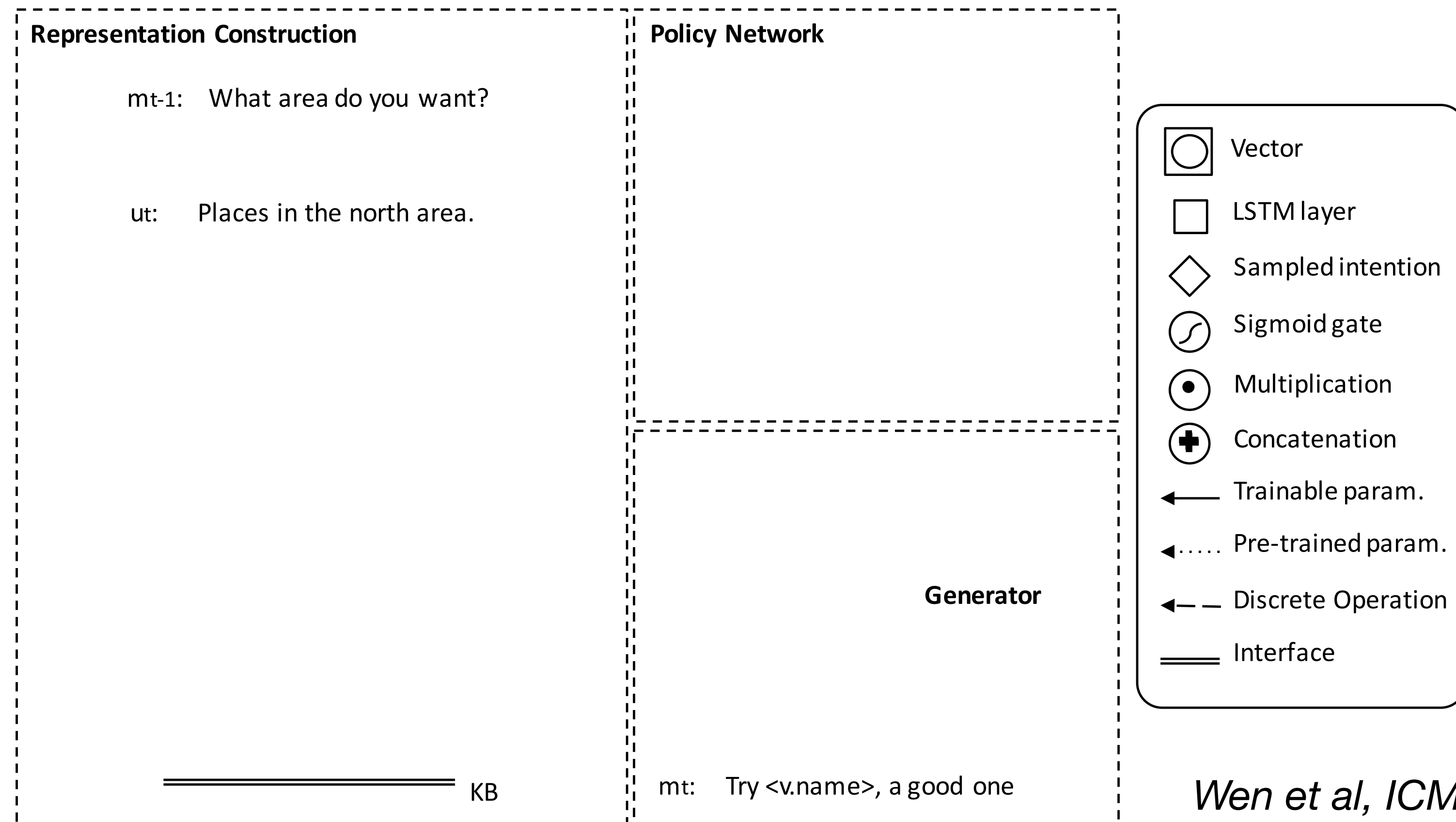- Task success rate is only around 41.1% when trained on 2000 dialogues.

POLYAi

# Neural Dialogue Model



- Explicitly use the structure of dialogue systems.

- Adopt slot-value pairs and domain-specific ontology in the development.

- A discrete access to database based on the parsed intent.

- Apply neural networks to predict intent and generate responses.

- Task success rate is around 92% when trained on only 400 dialogues.
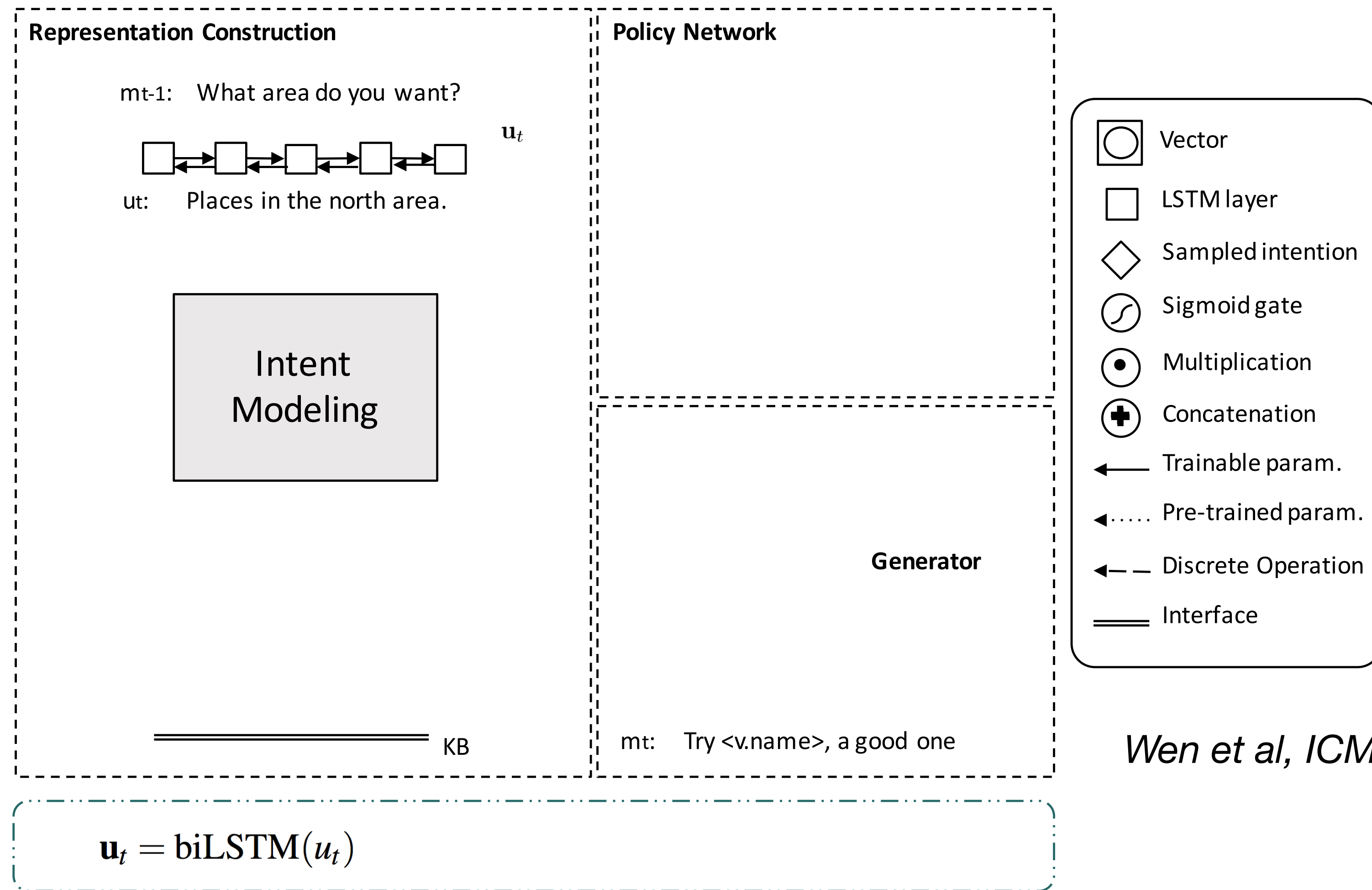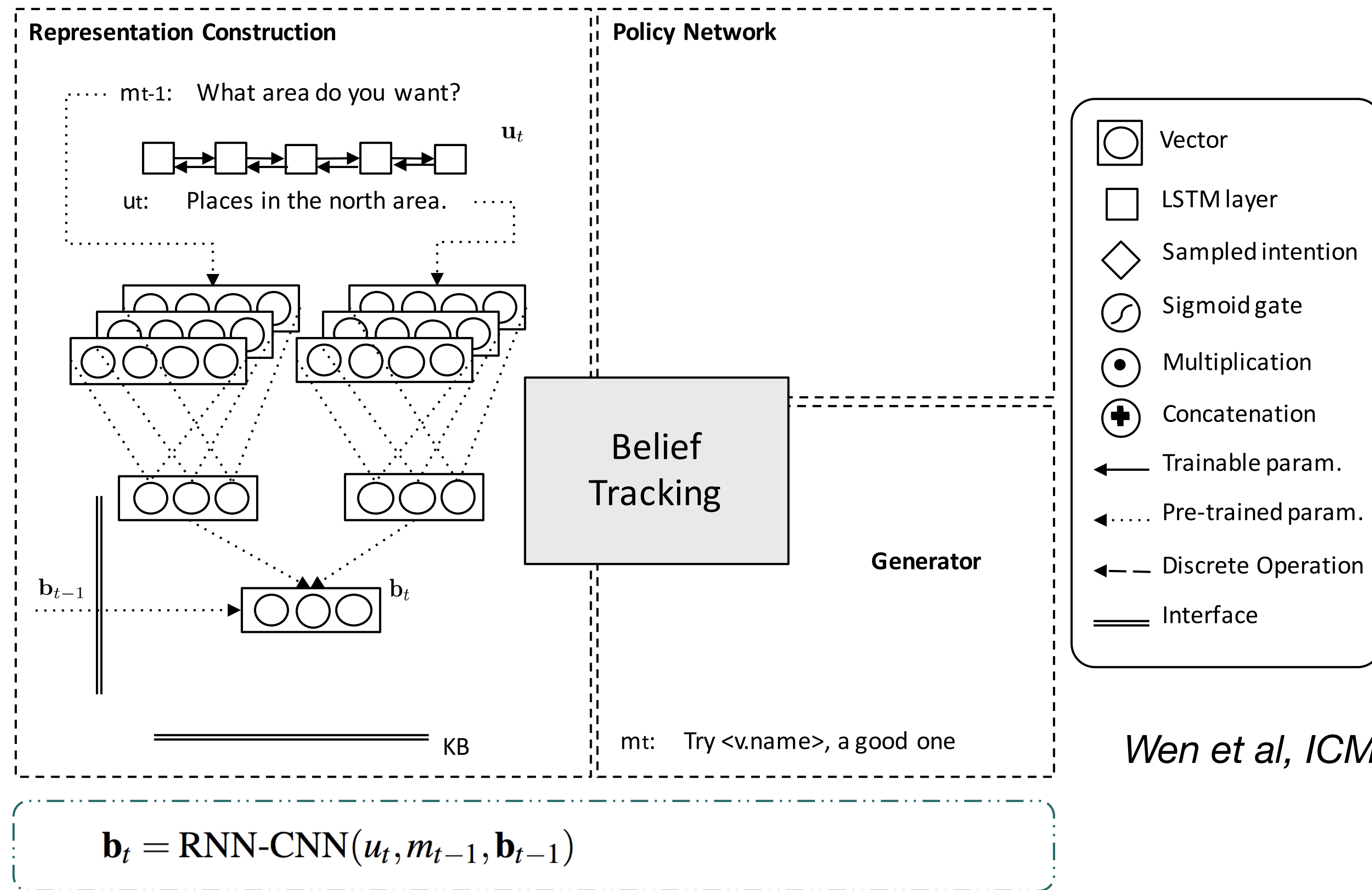
*Wen et al, EACL 2017*

# Latent Intention Dialogue Model

**Representation Construction**

m<sub>t-1</sub>:    What area do you want?

u<sub>t</sub>:    Places in the north area.

KB

**Policy Network**

**Generator**

m<sub>t</sub>:    Try <v.name>, a good one

Vector

LSTM layer

Sampled intention

Sigmoid gate

Multiplication

Concatenation

Trainable param.

Pre-trained param.

Discrete Operation

Interface

*Wen et al, ICML 2017.*

POLYAi

# Latent Intention Dialogue Model

**Representation Construction**

m$_{t-1}$:    What area do you want?

$\mathbf{u}_t$

u$_t$:    Places in the north area.

Intent Modeling

KB

**Policy Network**

**Generator**

m$_t$:    Try <v.name>, a good one

Vector

LSTM layer

Sampled intention

Sigmoid gate

Multiplication

Concatenation

Trainable param.

Pre-trained param.

Discrete Operation

Interface

*Wen et al, ICML 2017.*

$$\mathbf{u}_t = \mathrm{biLSTM}(u_t)$$

POLYAi

# Latent Intention Dialogue Model



$$\mathbf{b}_t = \text{RNN-CNN}(u_t, m_{t-1}, \mathbf{b}_{t-1})$$

*Wen et al, ICML 2017.*

# Latent Intention Dialogue Model



Wen et al, ICML 2017.

# Latent Intention Dialogue Model



Wen et al, ICML 2017.

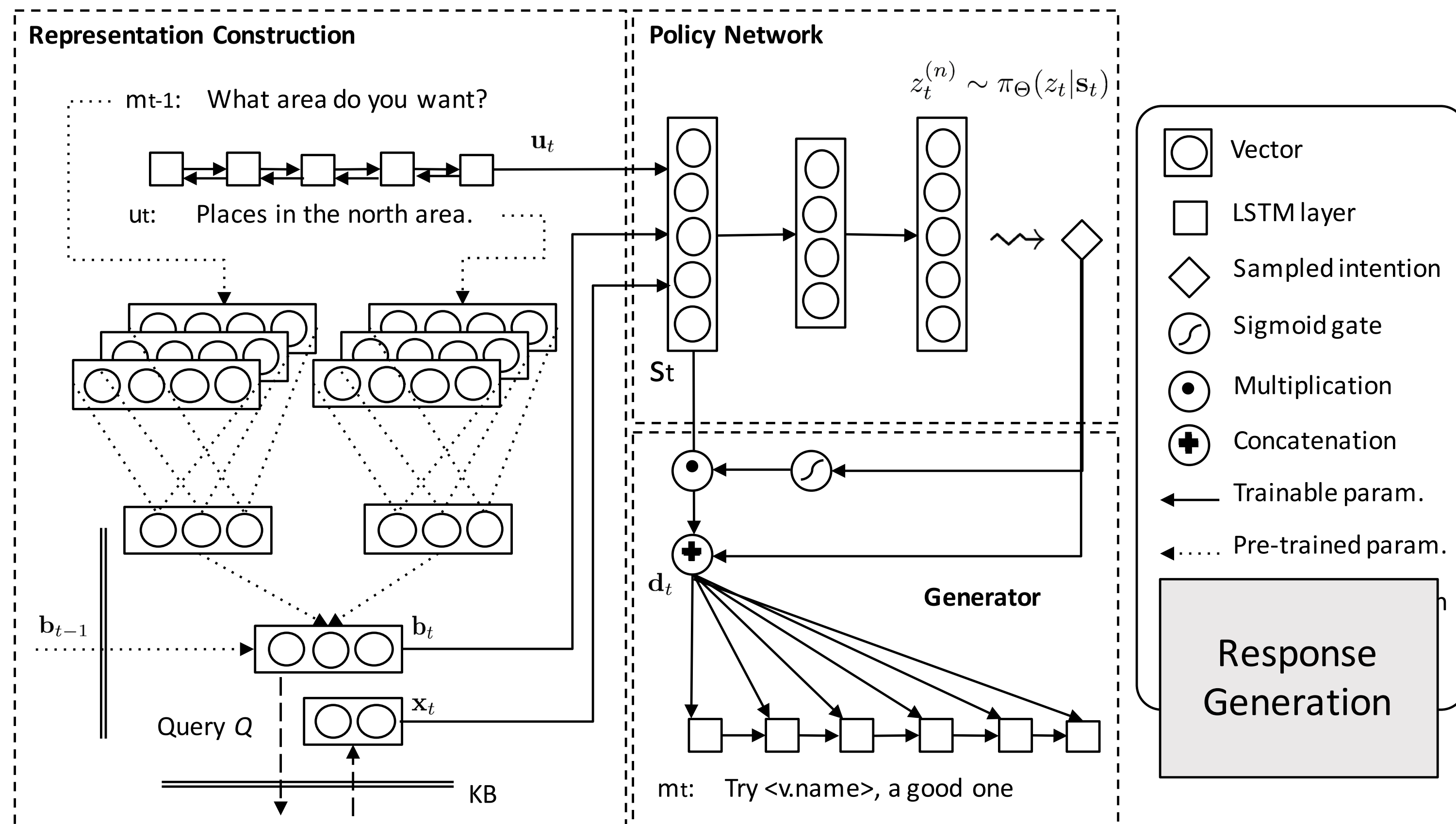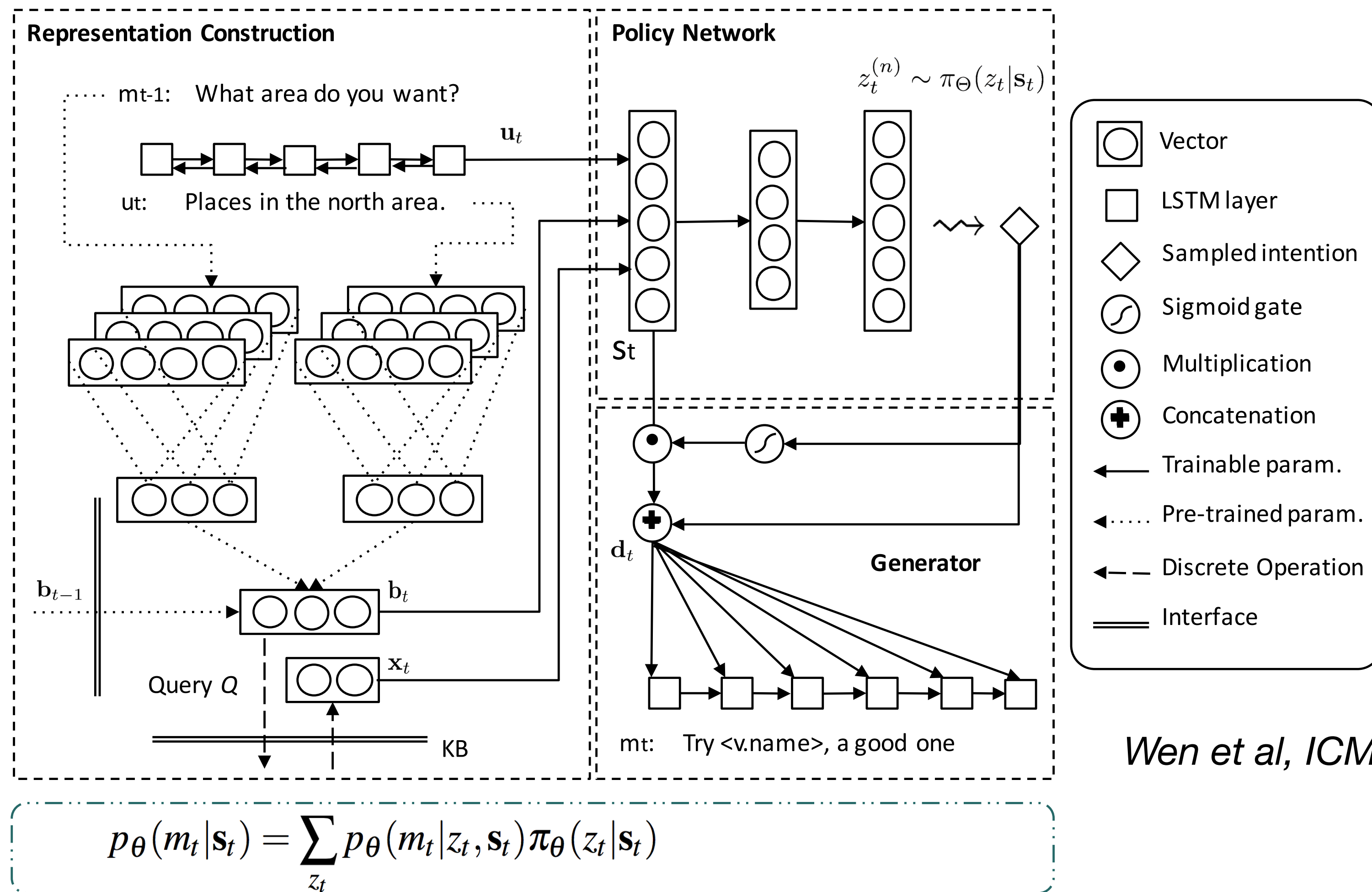$$\mathbf{s}_t = \mathbf{u}_t \oplus \mathbf{b}_t \oplus \mathbf{x}_t \qquad \pi_\theta(z_t | \mathbf{s}_t) = \mathrm{softmax}(\mathbf{W}_2^\mathsf{T} \cdot \tanh(\mathbf{W}_1^\mathsf{T}\mathbf{s}_t + \mathbf{b}_1) + \mathbf{b}_2)$$

# Latent Intention Dialogue Model



Representation Construction

mt-1:    What area do you want?

$\mathbf{u}_t$

ut:    Places in the north area.

$\mathbf{b}_{t-1}$    $\mathbf{b}_t$

Query $Q$    $\mathbf{x}_t$

KB

Policy Network

$z_t^{(n)} \sim \pi_\Theta(z_t | \mathbf{s}_t)$

St

Sampling as decision-making

Generator

mt:    Try <v.name>, a good one

Vector

LSTM layer

Sampled intention

Sigmoid gate

Multiplication

Concatenation

Trainable param.

Pre-trained param.

Discrete Operation

Interface

*Wen et al, ICML 2017.*

$$z_t^{(n)} \sim \pi_\theta(z_t | \mathbf{s}_t)$$

# Latent Intention Dialogue Model



$$\mathbf{d}_t = \mathbf{W}_4^\mathsf{T}\mathbf{z}_t \oplus \left[\mathrm{sigmoid}(\mathbf{W}_3^\mathsf{T}\mathbf{z}_t + \mathbf{b}_3) \cdot \mathbf{W}_5^\mathsf{T}\mathbf{s}_t\right] \qquad p_\theta(m_t|\mathbf{s}_t, z_t) = \prod_j p(w_{j+1}^t|w_j^t, \mathbf{h}_{j-1}^t, \mathbf{d}_t)$$

# Latent Intention Dialogue Model



Wen et al, ICML 2017.

$$p_\theta(m_t|\mathbf{s}_t) = \sum_{z_t} p_\theta(m_t|z_t, \mathbf{s}_t)\pi_\theta(z_t|\mathbf{s}_t)$$

# Learning of LIDM

- Approximate the objective with a variational lowerbound:

$$\mathcal{L}(\boldsymbol{\theta},\boldsymbol{\phi}) = \mathbb{E}_{q_\phi(z_t)}[\log p_\theta(m_t|z_t,s_t)] - \lambda D_{KL}(q_\phi(z_t)\|\pi_\theta(z_t|s_t))$$
$$\leq \log \sum_{z_t} p_\theta(m_t|z_t,s_t)\pi_\theta(z_t|s_t)$$
$$= \log p_\theta(m_t|s_t)$$

- Follow the **Variational Inference** method.

- *Latent Intention Dialogue Models, Wen et al, ICML 2017*

POLYAi

# Human Evaluation

| Metrics | NDM | LIDM | LIDM+RL |
|---|---|---|---|
| Success | 91.5% | 92.0% | 93.0% |
| Comprehension | 4.21 | $4.40^*$ | 4.40 |
| Naturalness | 4.08 | $4.29^*$ | $4.28^*$ |
| # of Turns | 4.45 | 4.54 | 4.29 |

\* $p < 0.05$

- A high dialogue success rate: >91%.

- Comprehension and naturalness are both rated high >4/5

POLYAi

# Sample Conversation

| Human : | Hi, I'm hungry for some Indonesian |
|---|---|
| Machine: | **( 0    0.71)    i am sorry , there are no indonesian restaurants in the area . would you like something else ?** |
| | (46    0.19)    there are no restaurants serving indonesian kind of food . |
| | (34    0.07)    what location would you like ? |
| Human : | Oh no. How about Indian? |
| Machine: | ( 0    0.59)    there are 9 restaurants in the east side of town . would you like their address ? |
| | (30    0.15)    curry prince serves indian cuisine . |
| | **(34    0.10)    what area would you like ?** |
| Human : | East side please |
| Machine: | **( 0    0.62)    curry prince is a indian restaurant in the east side of town . would you like their telephone number ?** |
| | (84    0.19)    curry prince is an indian restaurant in the east of town and the number is 01223 566388 . |
| | (81    0.08)    curry prince serves indian type of cuisine and is located in the east part of town . |
| Human : | yes , the address and phone number please |
| Machine: | ( 3    0.46)    the address is 451 newmarket road fen ditton and the phone is 01223 566388 . |
| | **(11    0.16)    the address for curry prince is 451 newmarket road fen ditton and the phone number is 01223 566388 .** |
| | ( 0    0.15)    curry prince is located at 451 newmarket road fen ditton and their phone is 01223 566388 . is there anything else i can help you with ? |
| Human : | Okay, thank you goodbye |
| Machine: | ( 0    0.57)    thank you , goodbye . |
| | ( 2    0.13)    goodbye . |
| | **( 1    0.11)    you are welcome . goodbye .** |
| | ( 6    0.10)    thank you for using the cambridge restaurant system . goodbye . |

*Table 4.* A sample dialogue from the LIDM, I=100 model, one exchange per block. Each latent intention is shown by a tuple (index, probability), followed by a decoded response. The sample dialogue was produced by following the responses highlighted in bold.
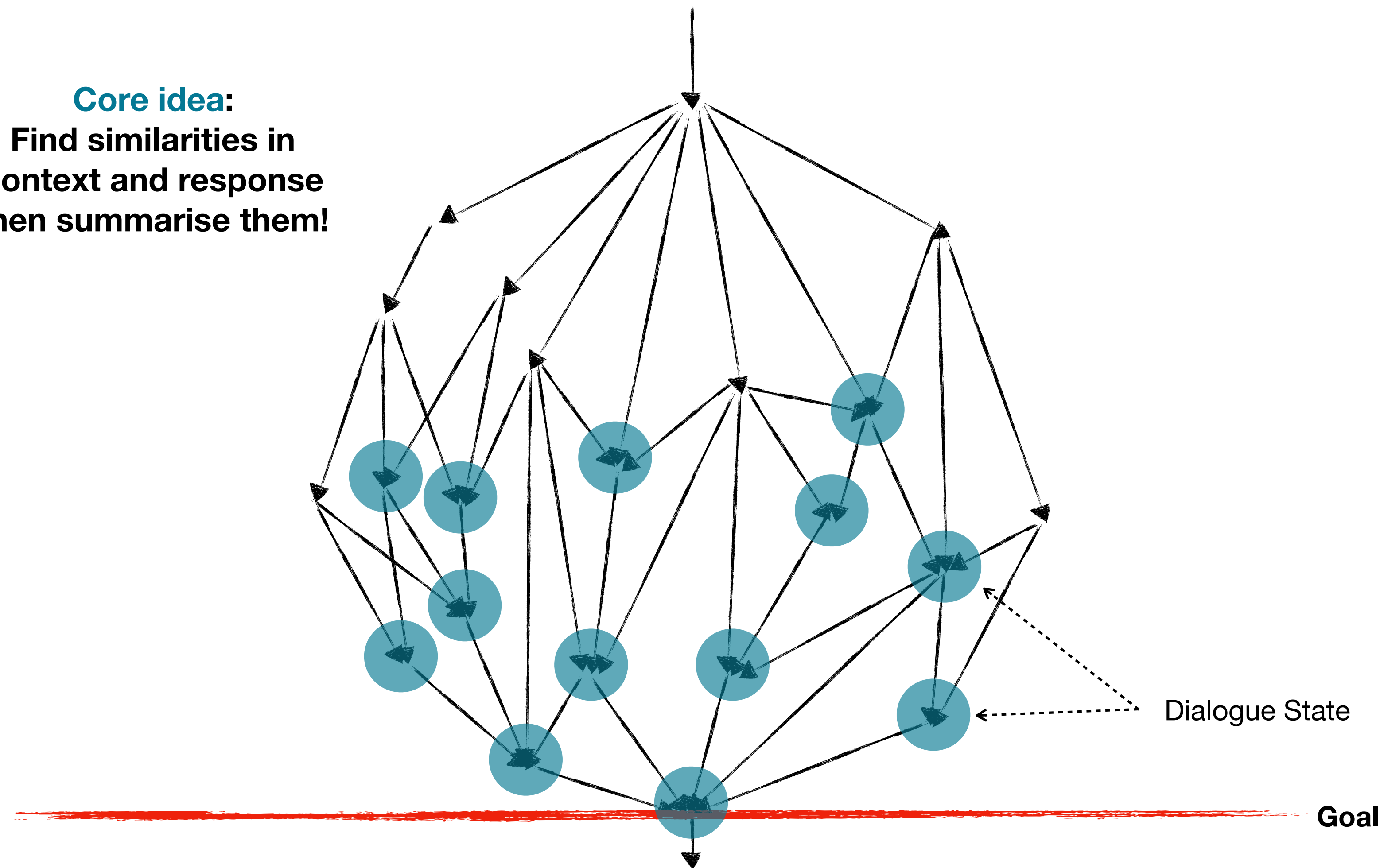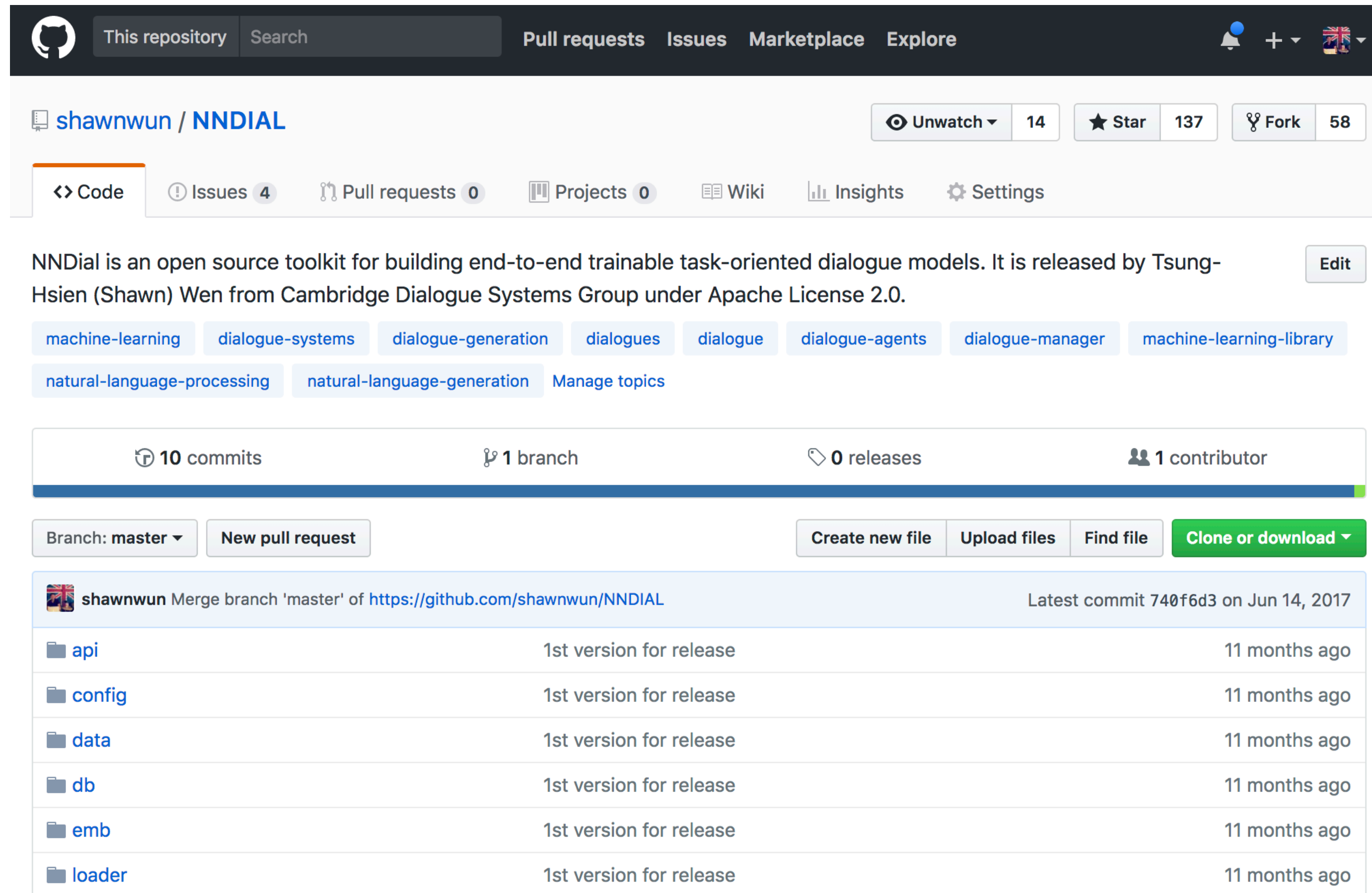
POLYAi

# Sample Conversation

| | | | |
|---|---|---|---|
| Human : | Hi, I'm hungry for some Indonesian | | |
| Machine: | **( 0** | **0.71)** | **i am sorry , there are no indonesian restaurants in the area . would you like something else ?** |
| | (46 | 0.19) | there are no restaurants serving indonesian kind of food . |
| | (34 | 0.07) | what location would you like ? |
| Human : | Oh no. How about Indian? | | |
| Machine: | ( 0 | 0.59) | there are 9 restaurants in the east side of town . would you like their address ? |
| | (30 | 0.15) | curry prince serves indian cuisine . |
| | **(34** | **0.10)** | **what area would you like ?** |
| Human : | East side please | | |
| Machine: | **( 0** | **0.62)** | **curry prince is a indian restaurant in the east side of town . would you like their telephone number ?** |
| | (84 | 0.19) | curry prince is an indian restaurant in the east of town and the number is 01223 566388 . |
| | (81 | 0.08) | curry prince serves indian type of cuisine and is located in the east part of town . |
| Human : | yes , the address and phone number please | | |
| Machine: | ( 3 | 0.46) | the address is 451 newmarket road fen ditton and the phone is 01223 566388 . |
| | **(11** | **0.16)** | **the address for curry prince is 451 newmarket road fen ditton and the phone number is 01223 566388 .** |
| | ( 0 | 0.15) | curry prince is located at 451 newmarket road fen ditton and their phone is 01223 566388 . is there anything else i can help you with ? |
| Human : | Okay, thank you goodbye | | |
| Machine: | ( 0 | 0.57) | thank you , goodbye . |
| | ( 2 | 0.13) | goodbye . |
| | **( 1** | **0.11)** | **you are welcome . goodbye .** |
| | ( 6 | 0.10) | thank you for using the cambridge restaurant system . goodbye . |

*Table 4.* A sample dialogue from the LIDM, I=100 model, one exchange per block. Each latent intention is shown by a tuple (index, probability), followed by a decoded response. The sample dialogue was produced by following the responses highlighted in bold.

POLYAi

# Sample Conversation

| | | | |
|---|---|---|---|
| Human : | Hi, I'm hungry for some Indonesian | | |
| Machine: | **( 0** | **0.71)** | **i am sorry , there are no indonesian restaurants in the area . would you like something else ?** |
| | (46 | 0.19) | there are no restaurants serving indonesian kind of food . |
| | (34 | 0.07) | what location would you like ? |
| Human : | Oh no. How about Indian? | | |
| Machine: | ( 0 | 0.59) | there are 9 restaurants in the east side of town . would you like their address ? |
| | (30 | 0.15) | curry prince serves indian cuisine . |
| | **(34** | **0.10)** | **what area would you like ?** |
| Human : | East side please | | |
| Machine: | **( 0** | **0.62)** | **curry prince is a indian restaurant in the east side of town . would you like their telephone number ?** |
| | (84 | 0.19) | curry prince is an indian restaurant in the east of town and the number is 01223 566388 . |
| | (81 | 0.08) | curry prince serves indian type of cuisine and is located in the east part of town . |
| Human : | yes , the address and phone number please | | |
| Machine: | ( 3 | 0.46) | the address is 451 newmarket road fen ditton and the phone is 01223 566388 . |
| | **(11** | **0.16)** | **the address for curry prince is 451 newmarket road fen ditton and the phone number is 01223 566388 .** |
| | ( 0 | 0.15) | curry prince is located at 451 newmarket road fen ditton and their phone is 01223 566388 . is there anything else i can help you with ? |
| Human : | Okay, thank you goodbye | | |
| Machine: | ( 0 | 0.57) | thank you , goodbye . |
| | ( 2 | 0.13) | goodbye . |
| | **( 1** | **0.11)** | **you are welcome . goodbye .** |
| | ( 6 | 0.10) | thank you for using the cambridge restaurant system . goodbye . |

*Table 4.* A sample dialogue from the LIDM, I=100 model, one exchange per block. Each latent intention is shown by a tuple (index, probability), followed by a decoded response. The sample dialogue was produced by following the responses highlighted in bold.

POLYAi

# The PolyAI Approach

**Core idea:**
Find similarities in context and response then summarise them!

Dialogue State

**Goal**

# NNDIAL Toolkit



**Link:**
https://github.com/shawnwun/NNDIAL

A toolkit for training neural
network based dialogue models.

Include both NDM, Attentive NDM,
and LIDM implementations.

Include a CamRest676 dialogue
dataset.

POLYAi

# RNNLG Toolkit



**Link:**
**https://github.com/shawnwun/RNNLG**

**A toolkit for rendering text from a dialogue act representation.**

**Include several RNN language generator implementation like Attention-based Encoder Decoder, SC-LSTM, H-LSTM … etc.**

**Include datasets from four different application domains.**

# Short Summary

- Conversational AI is a collective effort of multiple parties and requires a mix of multiple methods.

- Machine learning is

  - a perfect tool for understanding quries.

  - an OK method for response selection/generation.

  - but terrible for dialogue decision-making and database integration (the training data limitation).

- Practical but open-minded in choosing approaches is the way to build robust and truly useful conversational AIs of the future.

POLYAi

Data! Data! Data!

# The Importance of Data

- Data is the "*core*" of ML algorithms.

  - One of the reasons why Deep Learning took over.

- Data should be the "*core*" of an ML company too.

  - You should own your data.

  - You should have a way to generate your data.

  - You should have a way to label your data.

POLYAi

# The Importance of Data



AI Algorithm

~18 yrs

Good Data
Available

< 3 yrs

AI Breakthrough

**Human-Level Clean
Speech ASR (1994)**

**Human-level Control on
Atari Games (2015)**

**GoogLeNet at human-level
performance (2014)**

**Wall Street
Journal Corpus
(1991)**

**Arcade Learning
Environment
Dataset (2013)**

**ImageNet
Corpus (2010)**

**Hidden Markov
Model (1984)**

**Q-Learning
(1992)**

**Convolution
Network (1989)**

1980          1987          1994          2001          2008          2015

Reference: http://www.spacemachine.net/views/2016/3/datasets-over-algorithms

POLYAi

# The Importance of Data

- Conversational data is hard to get.

- Reddit has a large volume of conversational data.

  - For chit-chat dialogue development.

- But how do we get domain specific datasets?

  - A pizza delivery agent?

  - A taxi booking bot?

- COLLECT YOUR OWN DATA!

POLYAi

# Example: WOZ Data Collection
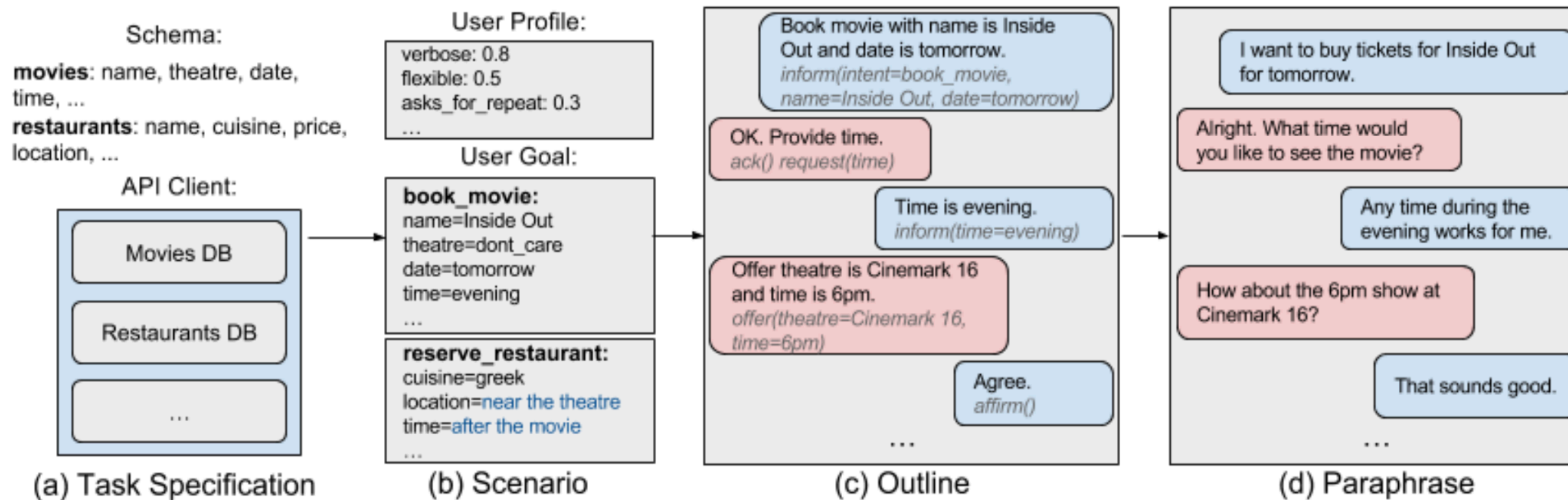


Wen et al, 2017

# Example: M2M Data Collection



Shah et al, 2018

# Short Summary

- Data is the lifeblood of machine learning algorithms.

- Historically, AI breakthroughs tend to correlate more to the datasets available in that area, rather than algorithms.

- In-domain conversational data is hard to get. Therefore, finding a way to collect/generate data is critical.

- PolyAI is building an in-house data collection platform to support client's business need.

POLYAi

# Conclusion

# Conclusion

- The dream of "*talking to machine*" is happening.

  - It is the matter how, not when.

- Existing tools do not provide the right kind of abstractions.

  - And there are so much to learn from the research community.

- A collective efforts are required to realize it

  - Across different industries.

  - Across multiple roles.

  - And multiple approaches.

POLYAi

POLYAi

# Q & A

Interested in **partnership** or **working** at PolyAI?
Send email to: contact@poly-ai.com

POLYAi