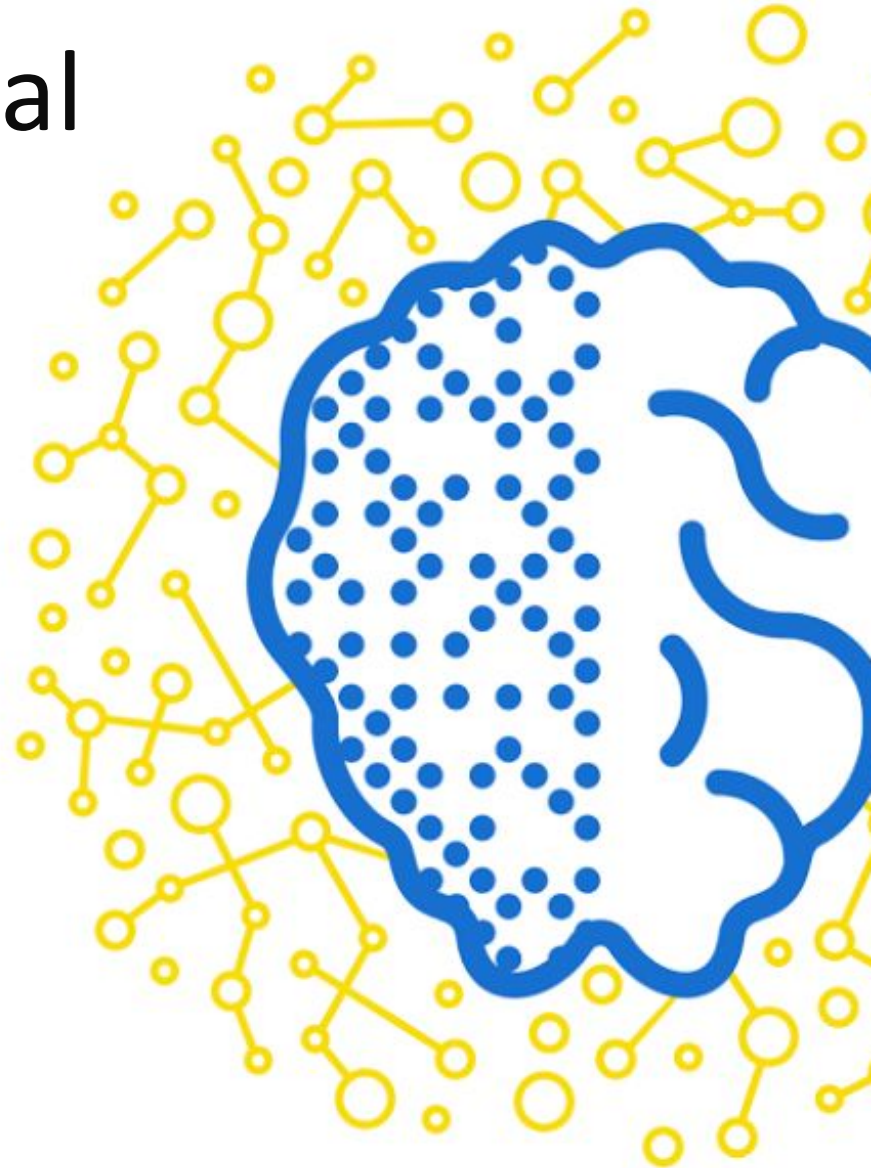# iPavlov: Conversational Intelligence Project

*Mikhail Burtsev, PhD*

*Moscow Institute of Physics and Technology (MIPT)*
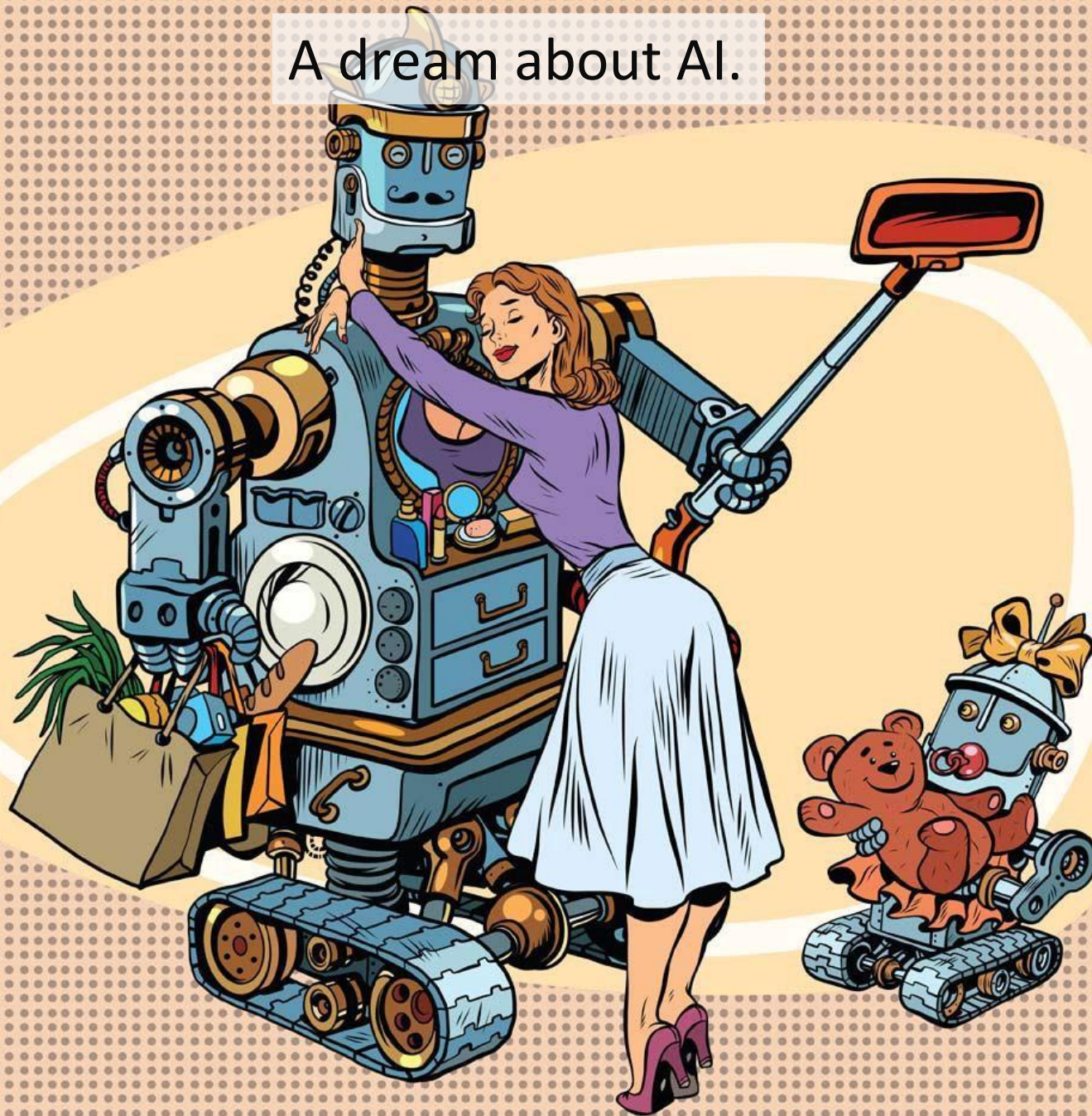
```python
# Definition of iPavlov project
def iPavlov(talent, ideas):
    research = ideas * talent
    AI = development(research)
    return AI
```
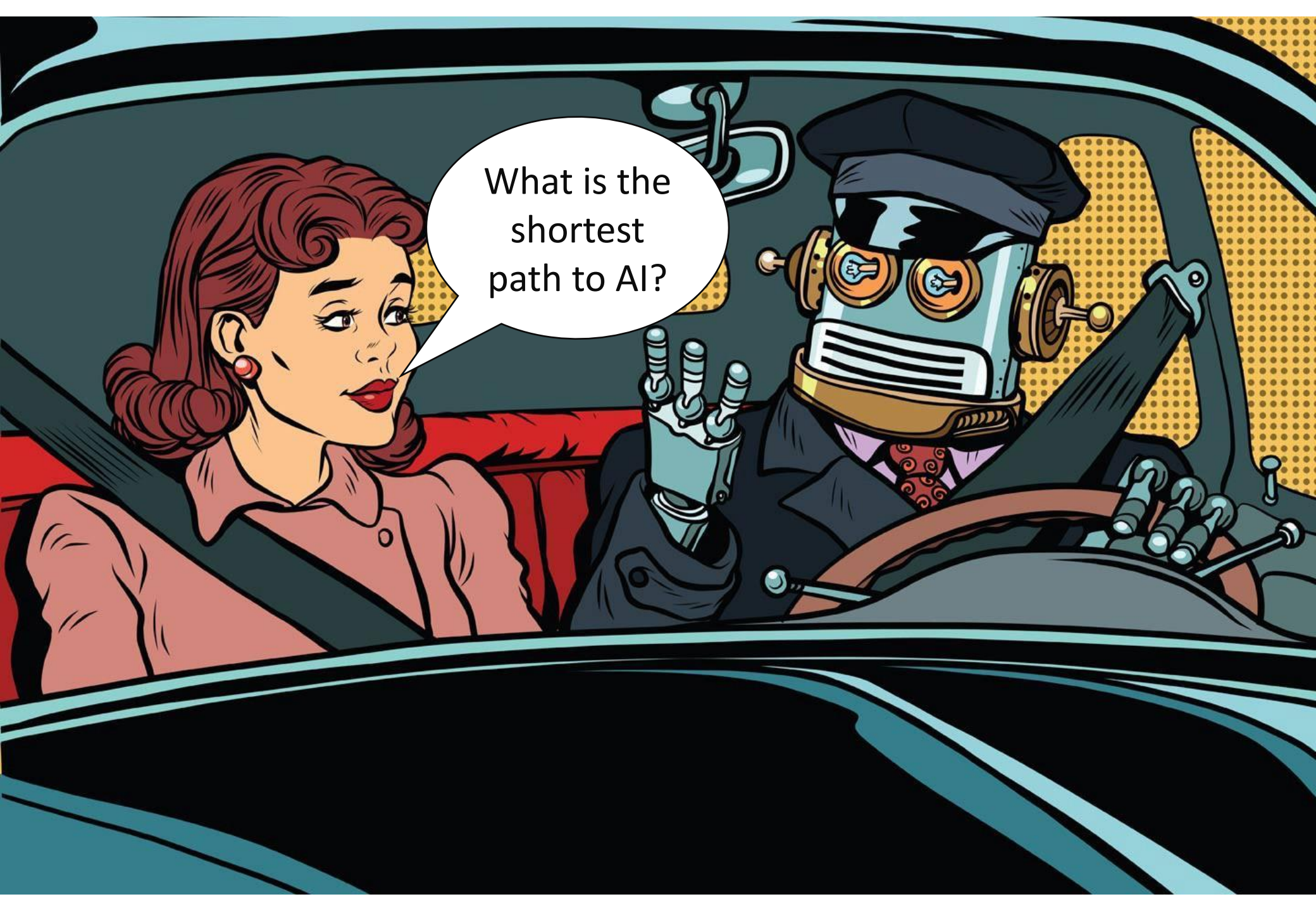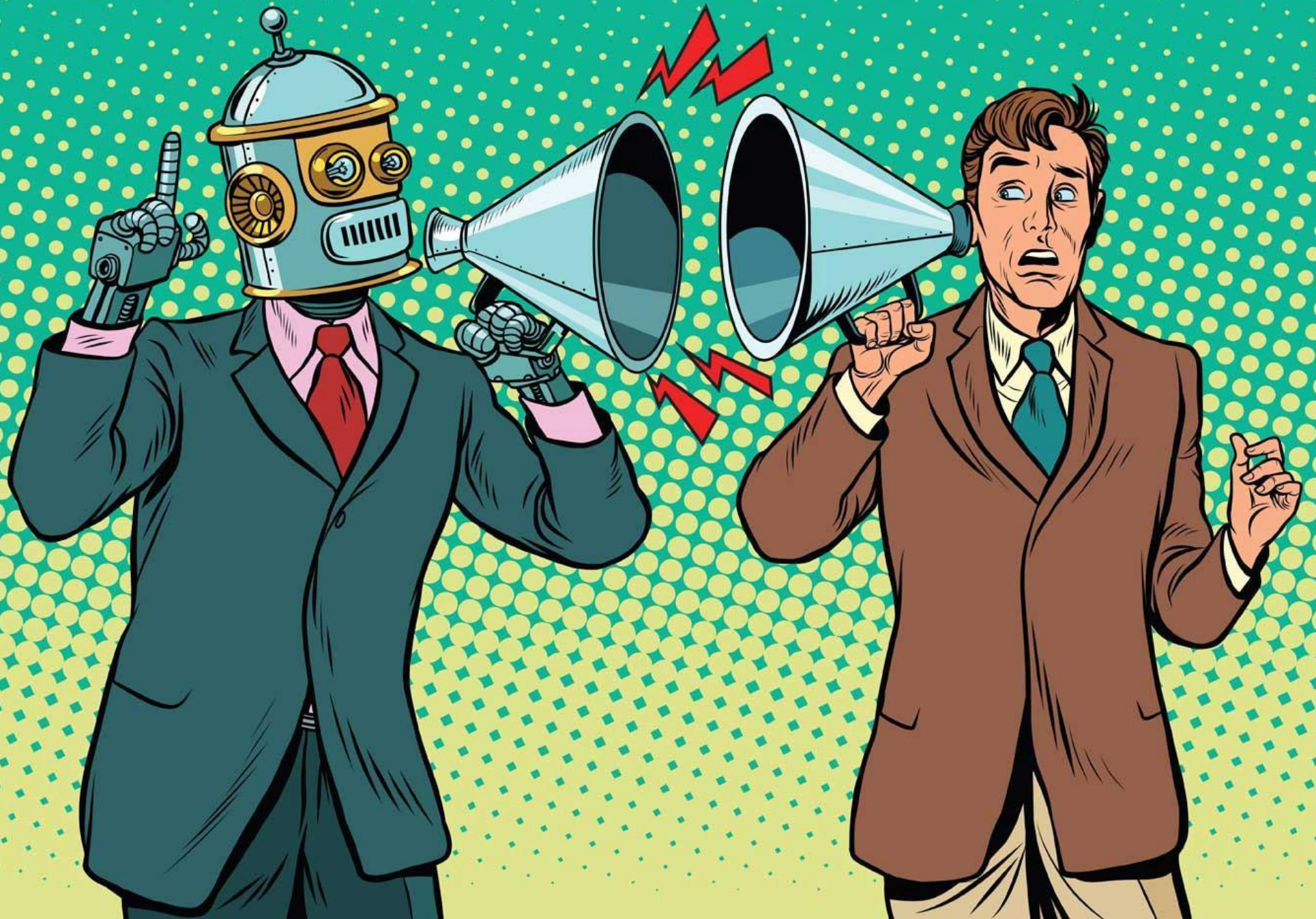
# Everybody has a dream

A dream about AI.

- Conversational Intelligence
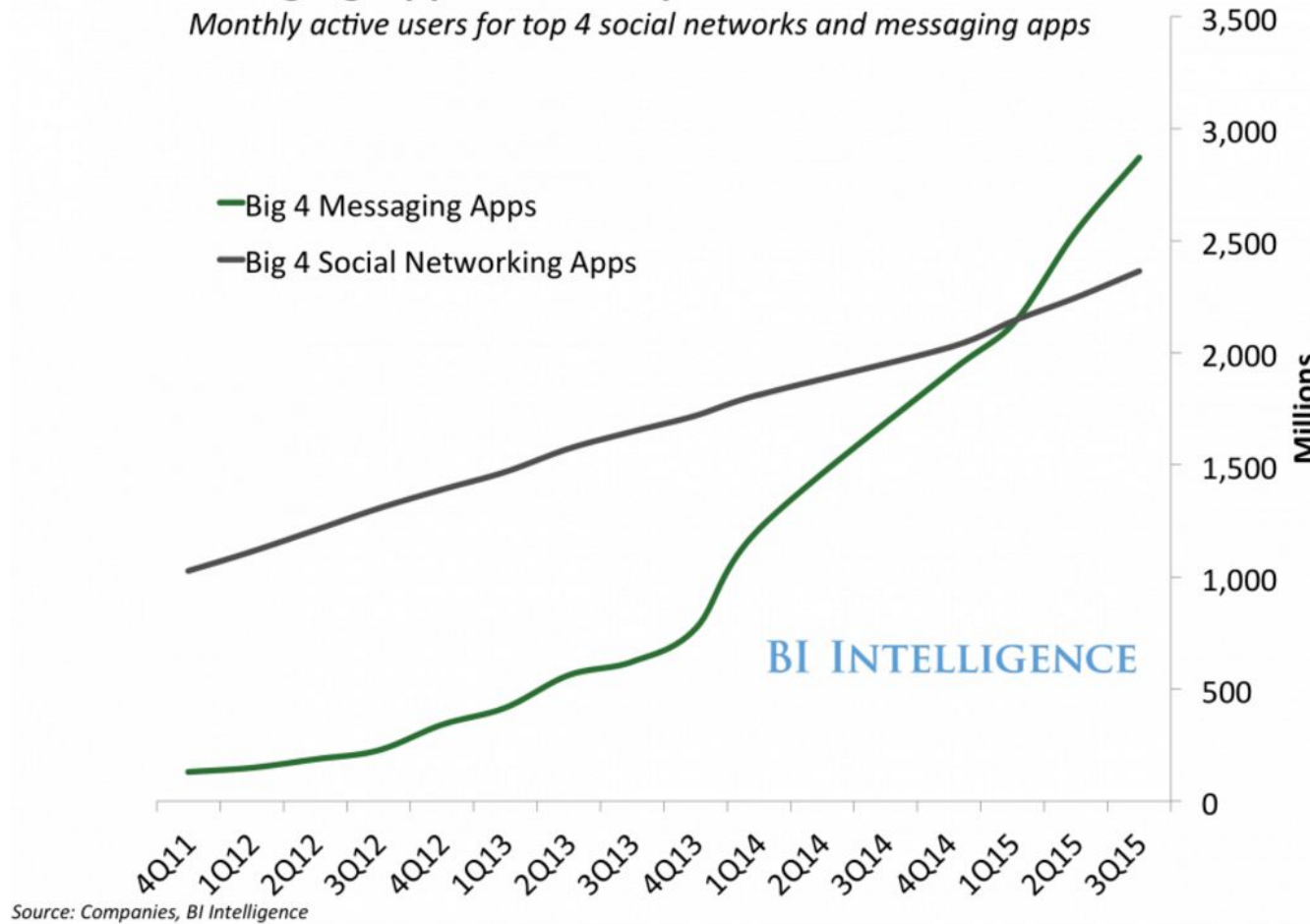
    - Complex real world problem

    - Can be decomposed into simpler tasks - NLU, DM, NLG

    - Big amount of data is available

    - Immediate application in industry

    - A step towards solving AI

- Promise of deep learning :

    - recurrent neural networks for the generation of sequences, and

    - attention and reinforcement learning for the dialogue planning .

iPavlov.ai

# Textual exchange dominates digital communication

## Messaging Apps Have Surpassed Social Networks

*Monthly active users for top 4 social networks and messaging apps*

— Big 4 Messaging Apps
— Big 4 Social Networking Apps

BI INTELLIGENCE

3,500
3,000
2,500
2,000
1,500
1,000
500
0

Millions

4Q11 1Q12 2Q12 3Q12 4Q12 1Q13 2Q13 3Q13 4Q13 1Q14 2Q14 3Q14 4Q14 1Q15 2Q15 3Q15

Source: Companies, BI Intelligence

# Conversational interface to seamlessly plug in human communication



Bots Landscape

DESIGNED BY JON CIFUENTES

POWERED BY VB | Profiles

iPavlov.ai

**iPavlov** project

*Deep learning architectures for the conversational intelligence*

- The major lab project for the 2017-2019

- Joint project with Sberbank the largest bank in Russia (operating income $20 billion, total assets $400 billion (2014))

- 20 researchers and engineers

**MIPT** MOSCOW INSTITUTE OF PHYSICS AND TECHNOLOGY

**SBERBANK** *By your side*

**Ivan Petrovich Pavlov** (1849 –1936) Russian physiologist known for his work in classical conditioning.

National Technølogy Initiative

Space of possibility

iPavlov.ai

**MIPT**
- AI Research Center

**SBERBANK**
*By your side*

**Sberbank**
- backend for AI powered applications

iPavlov

**Startup ecosystem**
- tools for rapid development of chat-bots

**Researchers**
- instruments for fast prototyping of models

iPavlov.ai

- ## Technology outcomes

  - Opensource deep learning NLP library **DeepPavlov**.

  - AI platform **DeepReply** implementing NLP services on top of **DeepPavlov** library for the chat-bot and dialogue systems products.

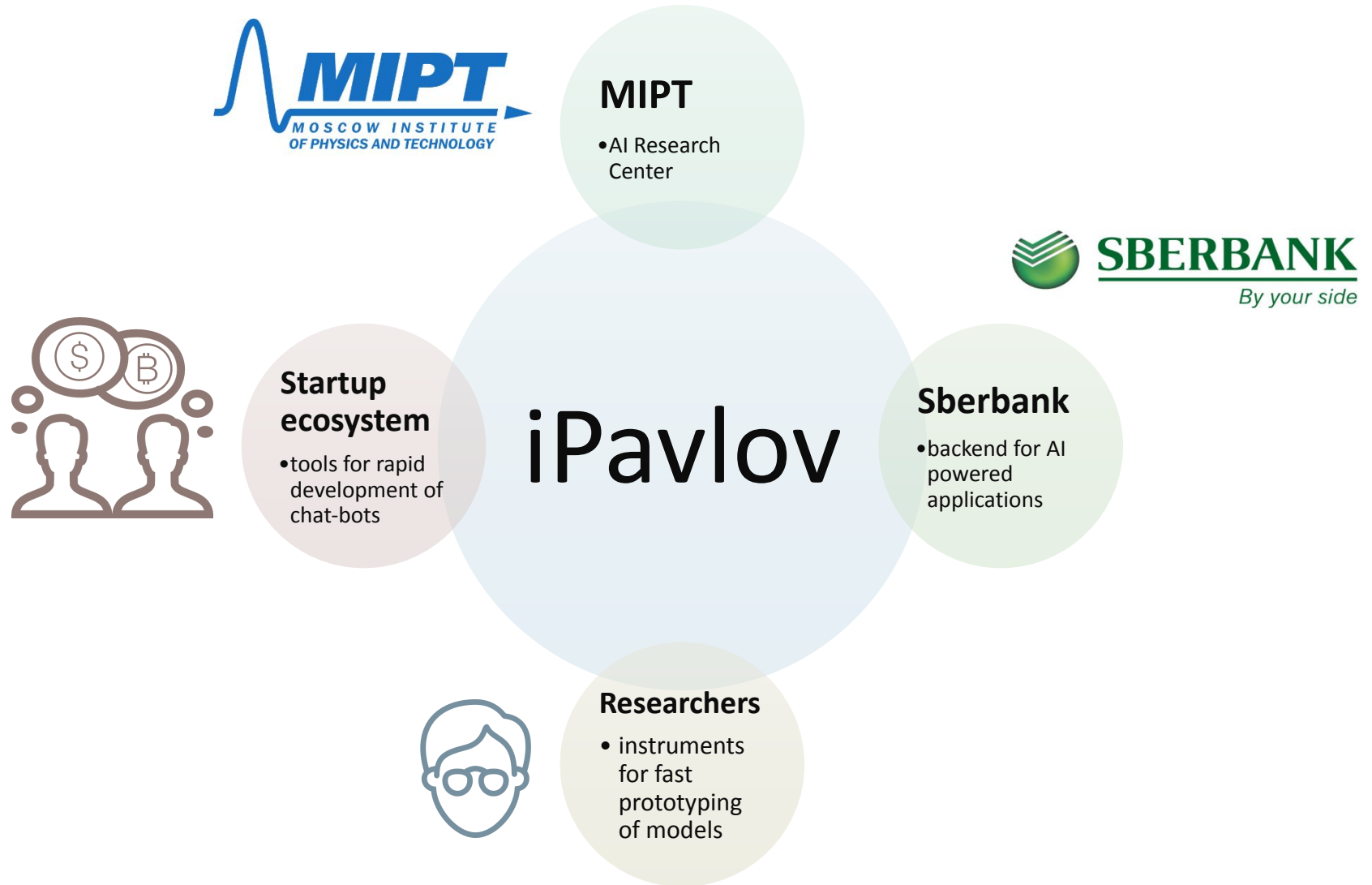| Technology Stack | Project Outcome | Description | Examples |
|---|---|---|---|
| AI APPLICATIONS | Out of the scope of iPavlov project | Third party AI applications in the domain of conversational intelligence. | Google Now, Digital Genius |
| AI SERVICES | **DeepReply** | AI conversational services to the neural network models trained for specific domains. | API.ai, wit.ai, Google NLP API |
| DEEP LEARNING ARCHITECTURES | **DeepPavlov** | Core components for neural conversational intelligence. Basic NLP functions and major neuroarchitectures for the dialogue systems. | MemNN, HRED |
| CORE DEEP LEARNING ALGORITHMS | | | Seq2seq, CNN, RNN, LSTM |
| COMPUTATIONAL LIBRARIES | Out of the scope of iPavlov project | | ThensorFlow (Google), Torch(Facebook), |
| DRIVERS GPU/FPGA | | | C/C++,Python, Julia… |
| CPU/GPU/FPGA | | | NVIDIA GPU, Intel CPU, Google TPU |

iPavlov.ai

## Research

Neural architectures for dialogue systems

Neural networks and reinforcement learning for planning

## Development
## **DeepPavlov**
## open source library

Repository of dialogue agents' models for variety of tasks

Lego-like modules for the fast prototyping of dialogue systems

Service NLP functions

## Applications
## **DeepReply**
## services

Conversational agents for specific business cases

API for separate NLU, DM, NLG tasks

# Modular dialog system



Are there any comedy movies to see this weekend?

*text data*

**NLU**
(Natural Language Understanding)
- Domain detection
- Intent detection
- Entities detection

intent = request_movie
entities = { genre = 'комедии',
            date = 'выходные '
}
*semantic frame*

Where are you?

*text data*

**NLG**
(Natural Language Generation)
- Generative models
- Templates

action = request_location

*system action*

**DM**
(Dialogue manager)
- Состояние диалога
- Политика поведения

iPavlov.ai

- ## Google Neural Machine Translation



Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation https://arxiv.org/abs/1609.08144 , Mon, 26 Sep 2016

iPavlov.ai

target chars: "e" "l" "l" "o"

output layer:
| 1.0 | 0.5 | 0.1 | 0.2 |
| 2.2 | 0.3 | 0.5 | -1.5 |
| -3.0 | -1.0 | 1.9 | -0.1 |
| 4.1 | 1.2 | -1.1 | 2.2 |

W_hy

hidden layer:
| 0.3 | 1.0 | 0.1 | -0.3 |
| -0.1 | 0.3 | -0.5 | 0.9 |
| 0.9 | 0.1 | -0.3 | 0.7 |

W_hh

input layer:
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 |

W_xh

input chars: "h" "e" "l" "l"

**Нейросеть 1**
2 слоя по 4096 ячеек GRU
Глубина back propagation = 50 шагов
Сеть обучалась 7 эпох
Объем обучающей выборки 2.5M слов субтитров

**Нейросеть 2**
2 слоя по 4096 ячеек GRU
Глубина back propagation = 100 шагов
Сеть обучалась 6 эпох
Объем обучающей выборки 11.06M слов субтитров

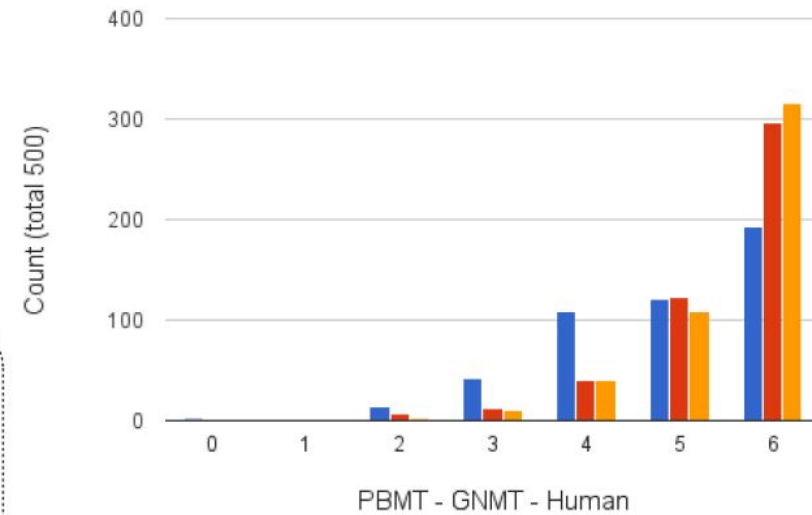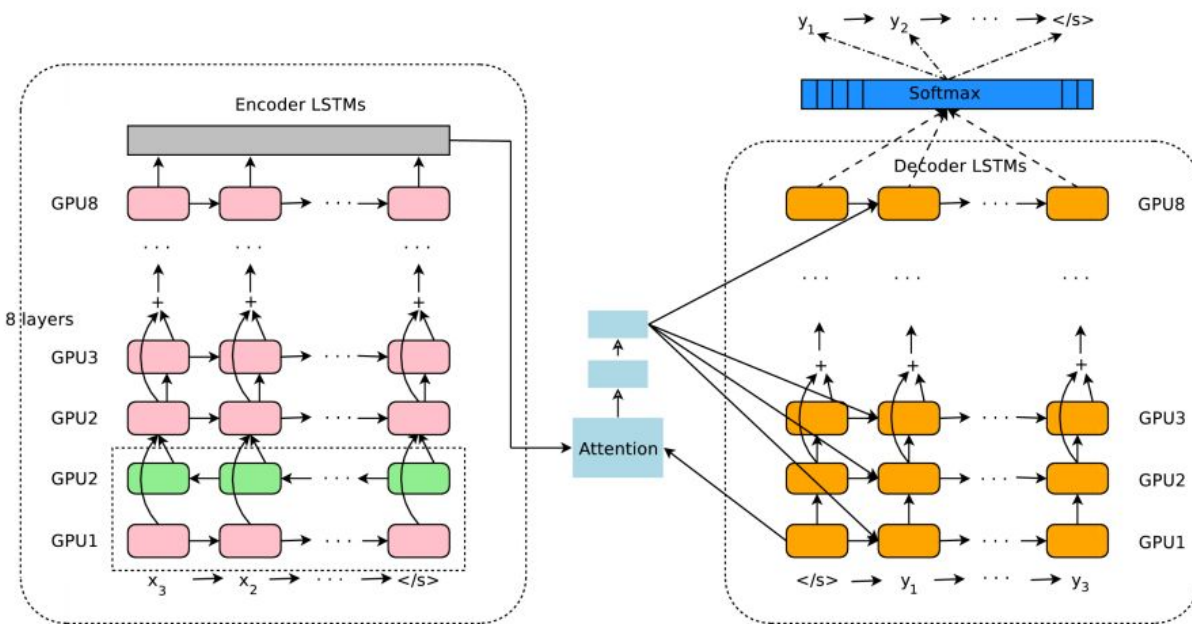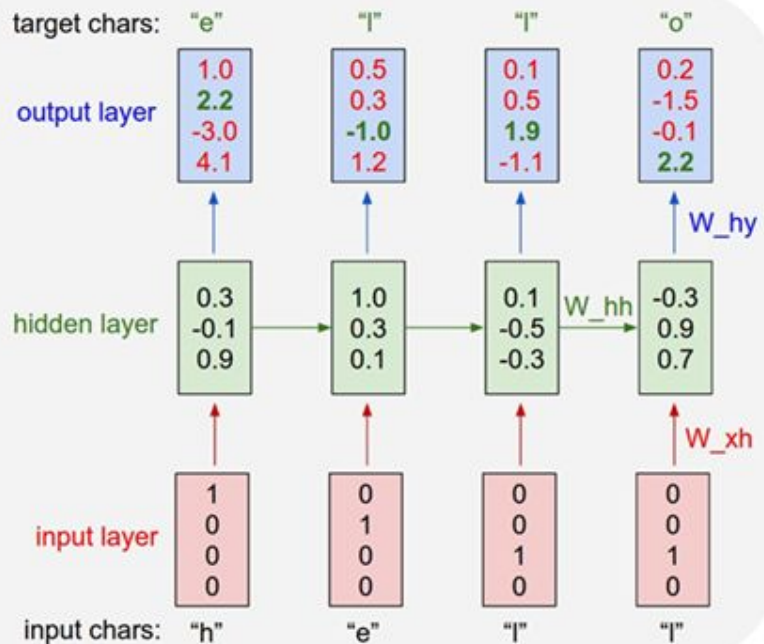Вопрос: Как ты устроен?
Ответ: Наверное, непросто.

Вопрос: Когда будет создан искусственный интеллект?
Ответ: Не знаю. Не особо важно.

Человек: Как построить искусственный интеллект?
Машина: Кто знает?

Человек: Первый закон робототехники
Машина: Удачного дня.

Человек: Ты глупый
Машина: Как насчет «привет».

Человек: Пошли гулять?
Машина: Это все так серьезно.
Человек: А что несерьёзно?
Машина: Что в этом плохого?
Человек: Ничего.
Машина: Кажется, наш парень пригласил меня.
Человек: Да, пригласил.
Машина: Не в этом дело.
Человек: А в чём?
Машина: Я просто хочу, чтобы вы все были снова в курсе всего, что ты делаешь.

Человек: Ты дура.
Машина: так почему тебе не все равно?
Человек: Потому что мне с тобой жить!
Машина: Ты совсем не такой!
Человек: Ты обиделась?
Машина: Угу.

Человек: Какой твой любимый фильм?
Машина: Это все личное дело.

Человек: Где найти нормального мужика?
Машина: Как ты себя чувствуешь?

Человек: Ты умеешь предсказывать будущее?
Машина: Что вы имеете в виду?
Человек: Можешь сказать, что случится завтра или через год?
Машина: Нам нужно выбираться отсюда.

Человек: Когда случится конец света?
Машина: Что ты имеешь в виду?
Человек: Вопросом на вопрос не отвечают.
Машина: Я тоже так думал.

Человек: зачем я сижу и читаю твои бесполезные ответы?
Машина: Я могу работать с тобой.

Человек: Ты веришь в бога?
Машина: Не знаю, как это может быть.

Человек: Забавно, но ты всего лишь программа и никогда не сможешь мыслить.
Машина: Конечно, дело в том, что сейчас не лучшее время.

Вопрос: Какой чудесный вечер. Не правда ли?
Ответ: Ты действительно думаешь, что я буду помнить об этом?

Человек: Как дела?
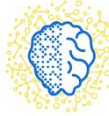Машина: Нам нужно их выдумать.

# Evolution of Neuro NLP Architectures

## Instant Encoder-Decoder

* Encode (read) one token with encoder network **E** then decode (write) one token with decoder network **D**.



$$\mathbf{X}_N^0 \quad \xrightarrow{\mathbf{x}_i} \quad E \quad \xrightarrow{\mathbf{e}_i} \quad \xrightarrow{\mathbf{e}_i} \quad D \quad \xrightarrow{\mathbf{d}_i} \quad \mathbf{Y}_N^0$$

Graves, A. (2013). Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850.*

## Encoder-Decoder (Seq2Seq)

* Encode (read) the whole sequence of tokens then decode (write) the whole sequence of tokens.
* Memory about the whole input sequence is encoded in the final state of the encoder **E**.



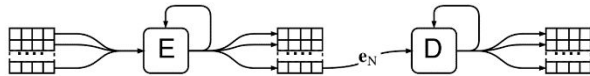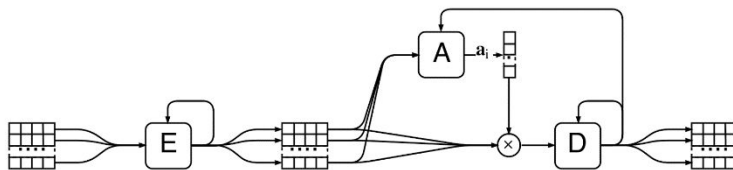Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078.

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in neural information processing systems* (pp. 3104-3112).

## Encoder-Decoder with Attention

* Encode (read) the whole sequence of tokens then decode (write) the whole sequence of tokens.
* Memory about every token of the input sequence is encoded and stored in a buffer separately.
* Attention sub-network **A** individually re-scales encodings of every input token taking into account the state of the decoder **D**.



Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473.*
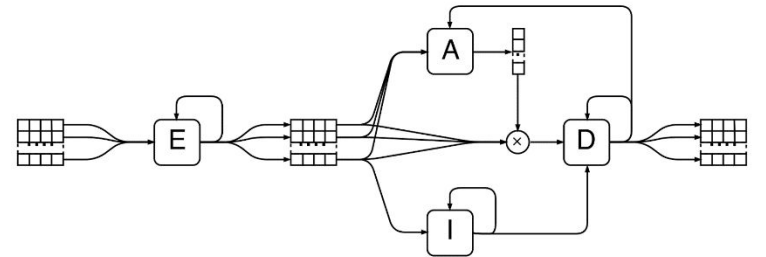
## Hierarchical Recurrent Encoder-Decoder (HRED)

* Encode (read) the whole sequence of tokens with encoder **E**$^1$ then update context memory **E**$^2$.
* Decode (write) the whole sequence of tokens with **E**$^2$ state as an additional input.



Serban, I. V., Sordoni, A., Bengio, Y., Courville, A., & Pineau, J. (2015). Building end-to-end dialogue systems using generative hierarchical neural network models. *arXiv preprint arXiv:1507.04808.*

## Attention with Intention Encoder-Decoder

* Encode (read) the whole sequence of tokens with **E**$^1$ then update context memory (intention) **I** (**E**$^2$).
* Decode (write) the whole sequence of tokens starting with **I** (**E**$^2$) state as initial hidden state of the decoder **D**.
* Attention sub-network **A** individually re-scales encodings of every input token taking into account the state of the decoder **D**.
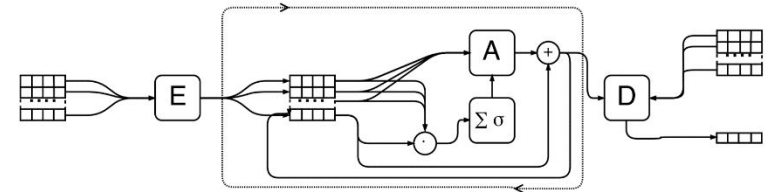


Yao, Kaisheng, Geoffrey Zweig, and Baolin Peng. "Attention with Intention for a NeuralNetwork Conversation Model." arXivpreprint arXiv:1510.08565 (2015).
Yao, Kaisheng, etal. "An Attentional Neural Conversation Model with ImprovedSpecificity." arXiv preprintarXiv:1606.01292 (2016).

## Memory Network

* Input is embedded sentences (replicas).
* Encode (read) the whole sequence of sentences' representations with linear encoding embedding **E** into memory.
* Encoding of the last sentence in memory is considered as "query" and controls "attention" **A**.
* Output of attention **A** is added to the old "query" to form a new query for the next iteration ("hop")
* After 1-3 iterations output of attention **A** is "compared" to possible candidate responses via linear "decoder" **D** and the best response is selected with softmax.



Bordes, A., & Weston, J. (2016). Learning end-to-end goal-oriented dialog. *arXiv preprint arXiv:1605.07683.*

iPavlov.ai

# Traditional pipeline in neural network implementation



**Natural Language Understanding**

**Dialog State Tracker**

**Policy**

**Natural Language Generation**

**Embedding or Encoder**:
mapping of input data to multidimensional space with desired properties resulting in vector representation

**Memory**:
history or context of the process represented as a set of vector representations

**Attention**:
given vector representation of the current input and memory controls hidden state of the system

**Decoder** or **Action generator**:
given hidden state of the system generates output

iPavlov.ai

# Sketch of the integrated architecture

- A year ago

iPavlov.ai

# Sketch of the integrated architecture



Memory Networks (Weston et.al., 2015)

HRED (Serban et.al., 2016)

long-term episodic memory

sentence encoder

sentence decoder

classifier

char level word encoder

char level word decoder

slot filling

QUERY

REPLY

iPavlov.ai

# Sketch of the integrated architecture



**Fine online learning**

episodic memory of sequences → generalization of episodic memory

episodic memory controller → evaluation of plans

**World model (planning)**

planning

dialog encoder → dialog decoder

**Policy**

**User simulation**

reply prediction

reply evaluation

**Reply re-evaluation**

dialog state corrector [2.2]

reply selection [2.3] ← hypothesis generator

**Dialog State Tracker**

long-term episodic memory

sentence encoder → sentence decoder

**NLU**

char level word encoder → char level word decoder

**NLG**

QUERY    REPLY

iPavlov.ai
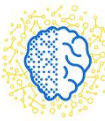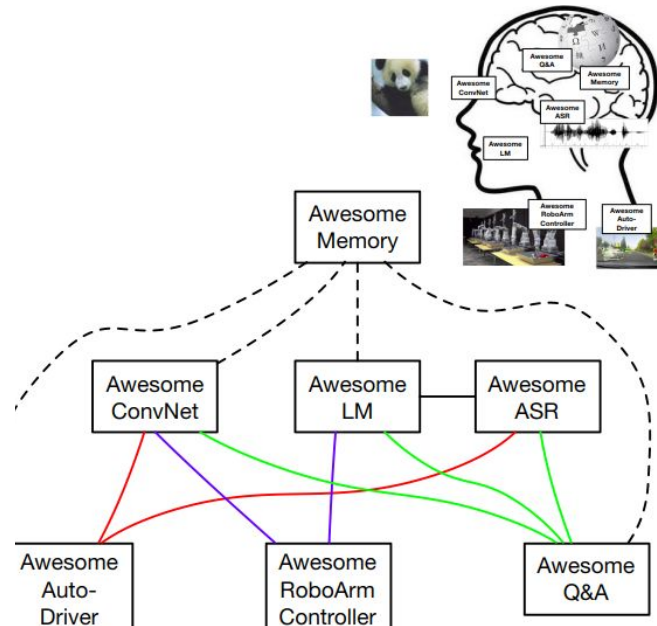
## What we want is...



- One system with many modules
- Modules interact with each other to solve a task
- Knowledge sharing across tasks via shared modules
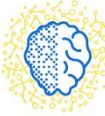- Some *trainable*, others *fixed*



## Paradigm shift

- One neural network *per task*

- One neural network *per function*
- Multiple networks cooperate to solve many higher-level tasks
- Mixture of trainable networks and fixed modules

Kyunghyun Cho (2017) *Deep Learning: a Next Step?*
https://drive.google.com/file/d/0B16RwCMQqrtdVWVGTE5LcWtwTzA/view

**iPavlov.ai**

# DeepPavlov

| Modules | Task-Oriented | Factoid | Chit-Chat |
|---|:---:|:---:|:---:|
| | T Agent | F Agent | C Agent |
| | **Task-Oriented** | **Factoid** | **Chit-Chat** |
| Named Entity Recognition | √ | √ | |
| Coreference resolution | √ | √ | |
| Paraphrase detection | √ | √ | |
| Insults detection | √ | | √ |
| Q&A | | √ | |
| Interactive Querying | √ | √ | |
| Memory | √ | | √ |
| Dialogue Policy | √ | | √ |
| … | | | |
| | DSTC-2 | SQuAD | reddit |

S Agent

iPavlov.ai

# DeepPavlov Open Source Library

iPavlov.ai

# Some results

- ## Named entity recognition in Russian



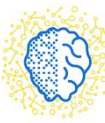| Models | Gareev's dataset | | | Persons-1000 | | | FactRuEval 2016 | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F |
| Gareev et al. [4] | 67.98 | 75.05 | 84.11 | - | - | - | - | - | - |
| Malykh et al. [9] | 59.65 | 65.70 | 62.49 | - | - | - | - | - | - |
| Trofimov [5] | - | - | - | 97.26 | 93.92 | 95.57 | - | - | - |
| Rubaylo et al. [21] | - | - | - | - | - | - | 77.70 | 78.50 | 78.13 |
| Sysoev et al.[8] | - | - | - | - | - | - | **88.19** | 64.75 | 74.67 |
| Ivanitsky et al. [7] | - | - | - | - | - | - | - | - | **87.88** |
| Mozharova et al. [6] | - | - | - | - | - | 97.21 | - | - | - |
| NeuroNER | 88.19 | 82.73 | 85.37 | 96.38 | 96.83 | 96.60 | 80.49 | 79.23 | 79.86 |
| NeuroNER + Highway char | 85.75 | **88.40** | 87.06 | 96.56 | 97.11 | 96.83 | 80.59 | 80.72 | 80.66 |
| NeuroNER + Highway LSTM | 84.35 | 81.96 | 83.14 | 96.49 | 97.19 | 96.84 | 81.09 | 79.31 | 80.19 |
| NeuroNER + Highway char + Highway LSTM | 83.33 | 85.05 | 84.18 | 96.74 | 96.83 | 96.78 | 79.13 | 78.76 | 78.95 |
| Bi-LSTM + CRF + *Lenta* | **89.57** | 84.89 | **87.17** | **99.43** | **99.09** | **99.26** | 83.88 | **80.84** | 82.10 |

Anh L., Arkhipov M., Burtsev M. Application of a Hybrid Bi-LSTM-CRF model to the task of Russian Named Entity Recognition // In proc. AINL, 2017
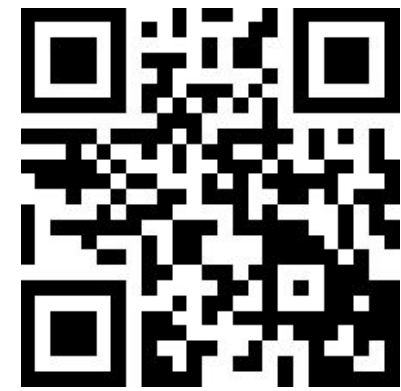
- Intent recognition

iPavlov.ai

# Challenges

- How to set goals in Task-Oriented neural end-to-end system?

- How to build a user model and integrate it with a dialogue agent?

- How to plan a dialogue with NN and RL implementation?

- How to evaluate dialogue systems?

- How to balance goal-directedness with engagement?

- How to integrate external information from DB, KB, IR un a dialogue?

- How to integrate modules and train integrated system?

- How to transfer knowledge from task to task?

- How to learn on-line?

iPavlov.ai

- Telegram @ConvaiBot

http://t.me/ConvaiBot

- Web page http://convai.io

- Dialog dataset http://convai.io/data/

## The Conversational Intelligence Challenge
### NIPS 2017 Live Competition

**Dialogue systems** and **conversational agents** –
including chatbots, personal assistants and voice control interfaces –
are becoming increasingly widespread in our daily lives.

**NIPS** is sponsoring an open competition to create a chatbot that can
hold an **intelligent conversation with a human partner**.

# Summary

- Textual user interface is becoming more and more intelligent

- Conversational intelligence evolves from modular towards end-to-end architectures

- iPavlov is R&D project with the goal to speed up prototyping of dialogue system for business and research

- DeepPavlov is an open source framework for the conversational intelligence

    - Repository of architectures for dialogue agents

    - Neural network components implementing NLU, DST, Policy, NLG and their combinations

- NIPS conversational challenge is an attempt to address the problem with dialogue systems evaluation

- Integration of IR and CI is the next step towards AI

iPavlov.ai

# iPavlov.ai



```python
# Definition of iPavlov project
def iPavlov(talent, ideas):
    research = ideas * talent
    AI = development(research)
    return AI
# How you are related to the iPavlov project
email.send('merge@ipavlov.ai', YOU.CV)
if YOU in ['researcher',
           'developer']
    and YOU is ('ai_geek' &
                'performer' &
                'team_player'):
        iPavlov(YOU.talent, YOU.ideas)
```

https://github.com/deepmipt/deeppavlov/