

# Сделаем экспериментально-опытное обучение интерпретируемым пожизненно!

*(Make experiential life-long learning interpretable!)*

Антон Колонин  
[akolonin@aigents.com](mailto:akolonin@aigents.com)  
Telegram: akolonin

**N\*** Novosibirsk  
State  
University  
\*THE REAL SCIENCE  
<https://www.nsu.ru>



<https://agirussia.org>

# Есть ли у нас проблемы?

Медленно обучаемся?

=> slow learning

Не можем объяснить свои решения?

=> uninterpretable models

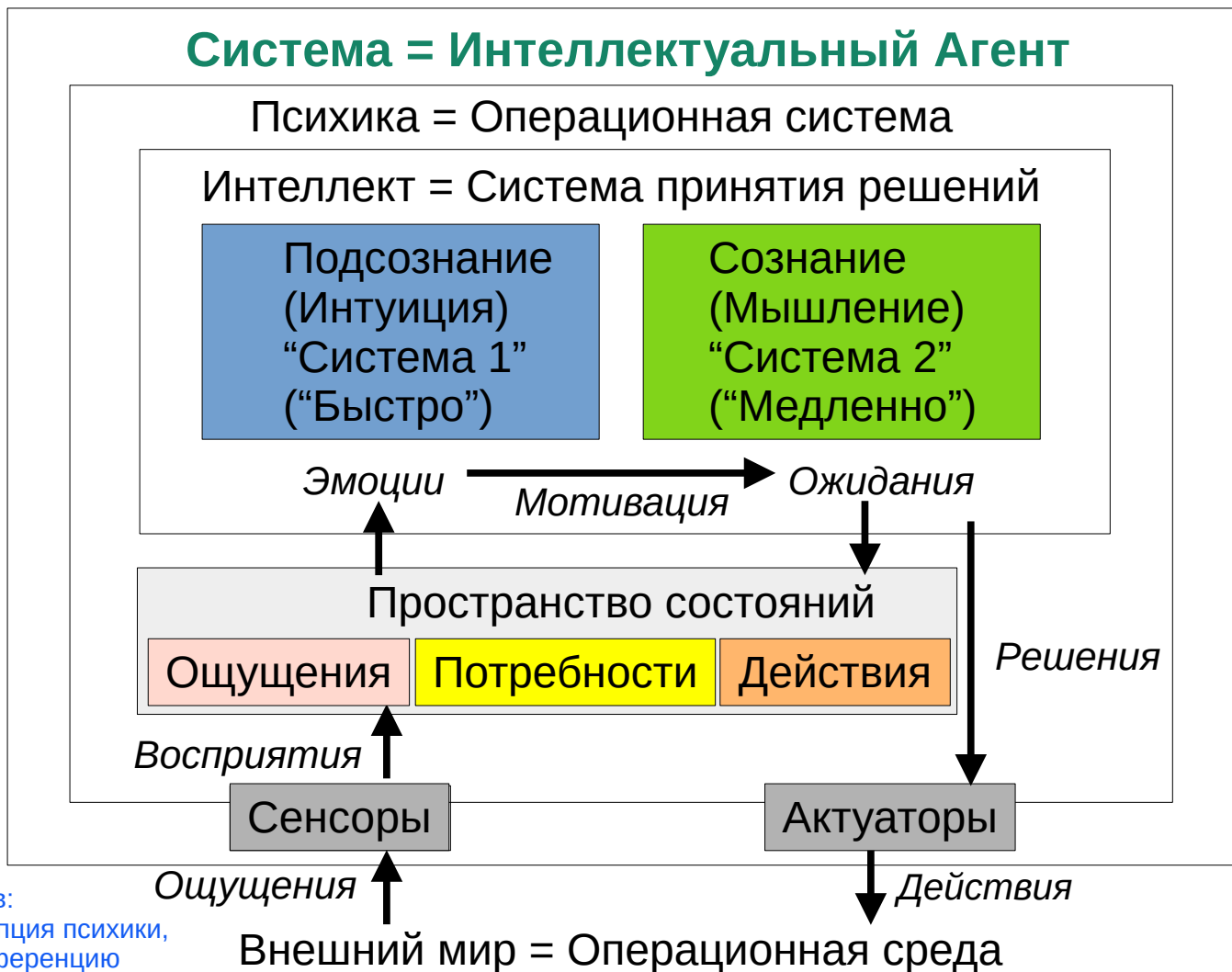
Забываем важное?

=> catastrophic forgetting

Тратим много энергии?

=> expensive, resource-consuming training

# Система = Интеллектуальный Агент

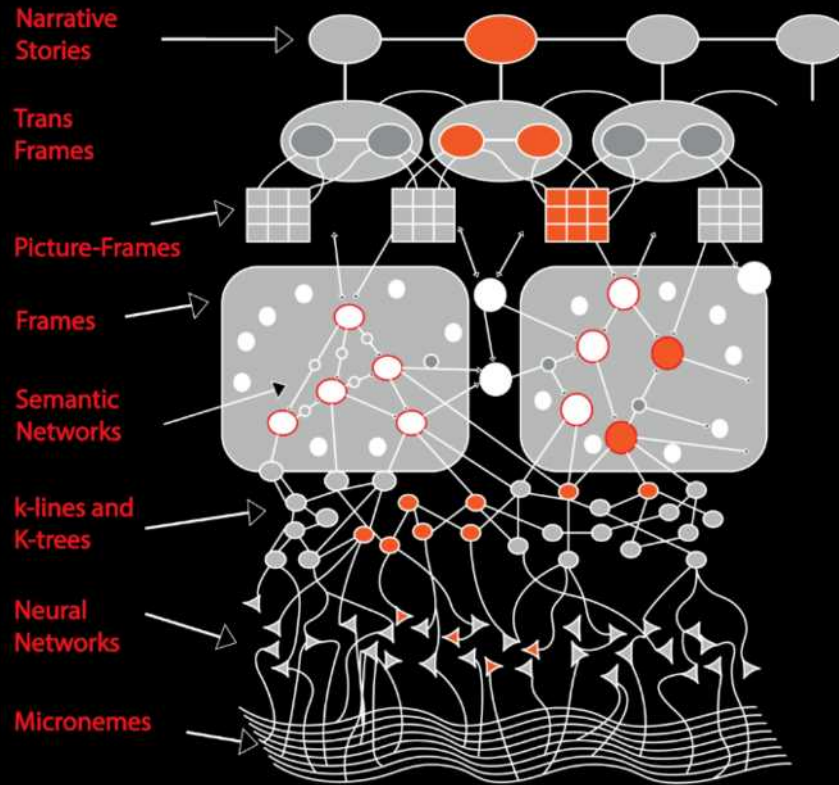


А.Г.Колонин, В.Г.Крюков:  
Вычислительная концепция психики,  
Статья принята на конференцию  
Нейроинформатика-25

# “Быстрое и медленное мышление” – Daniel Kahneman

easy  
explanation  
learning fast

hard  
explanation  
learning slow



thinking slow

thinking fast

<https://www.linkedin.com/pulse/explainable-ai-vs-explaining-part-1-ahmad-haj-mosa/>

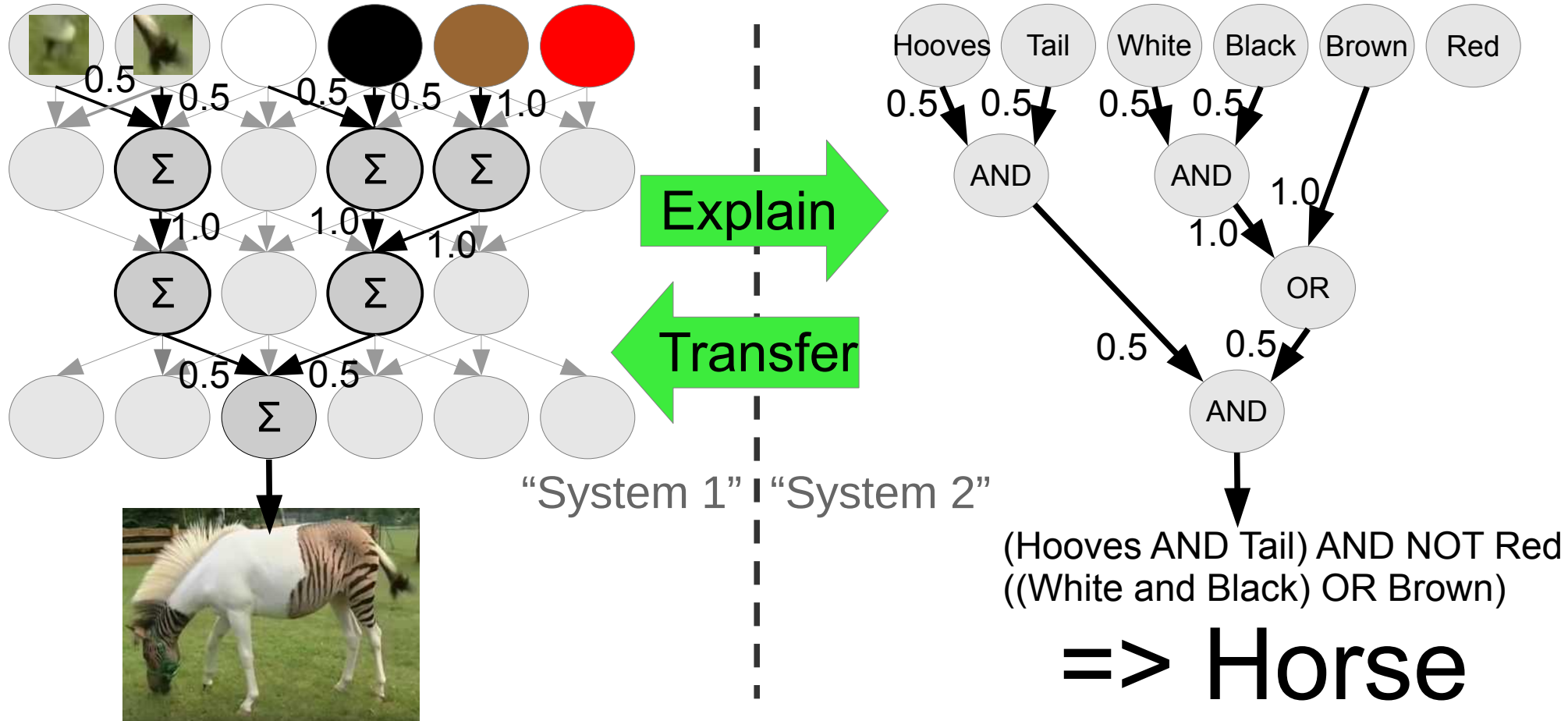
Xing, F., Cambria, E., Welsch, R. (2019). Theoretical Underpinnings on Text Mining. In: Intelligent Asset Management. Socio-Affective Computing, vol 9. Springer, Cham.

[https://doi.org/10.1007/978-3-030-30263-4\\_3](https://doi.org/10.1007/978-3-030-30263-4_3)

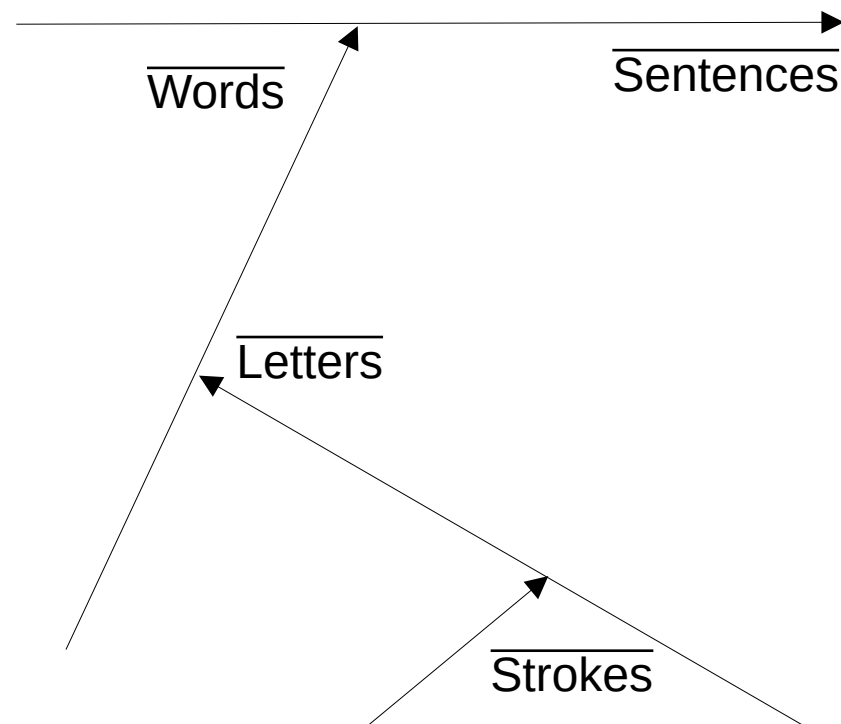
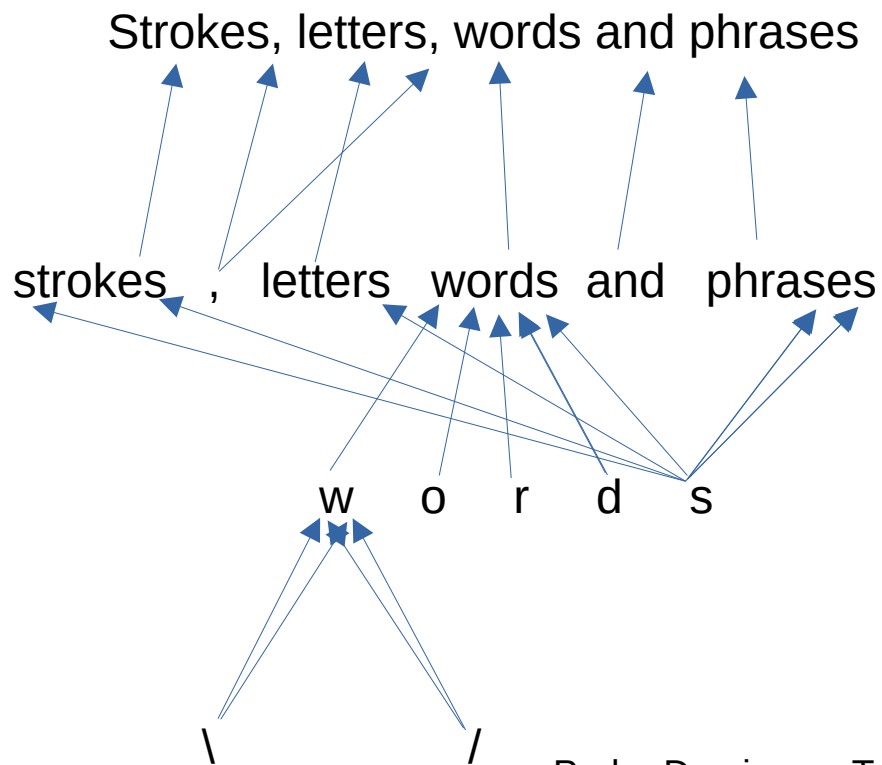
M. Minsky, The Emotion Machine: Commonsense Thinking, Artificial Intelligence, and the Future of the Human Mind (Simon & Schuster Paperbacks, Princeton, 2007)

Copyright © 2026 Anton Kolonin, Aigents®

# Нейро-символьная интеграция для интерпретируемого ИИ



# Функциональная эквивалентность графовых и ~~нейросетевых~~ тензорных моделей



Pedro Domingos, Tensor Logic: The Language of AI  
<https://arxiv.org/pdf/2510.12269>

Copyright © 2026 Anton Kolonin, Aigents®

Вершина графа - измерение  
Ребро графа - вектор  
Гиперграф - тензор 6

# Типизированная тензорная логика для различных систем ИИ (logical, sub-symbolic, probabilistic/non-axiomatic)

**Truth-Value Tensor**  
(NARS/PLN/...)

Property **0.0123456**  
**=750/60750**



striped  
horse  
Subject



Life-long  
learning?

**Numerical Tensor**  
(ANN/Bayesian Logic)

Property **~0.01**



striped  
horse  
Subject

**Boolean Tensor**  
(Boolean Logic)

Property **False**



striped  
horse  
Subject

Pei Wang: Non-Axiomatic Logic  
<https://www.worldscientific.com/worldscibooks/10.1142/14486>

Pedro Domingos, Tensor Logic: The Language of AI  
<https://arxiv.org/pdf/2510.12269>

# Психика = Операционная система

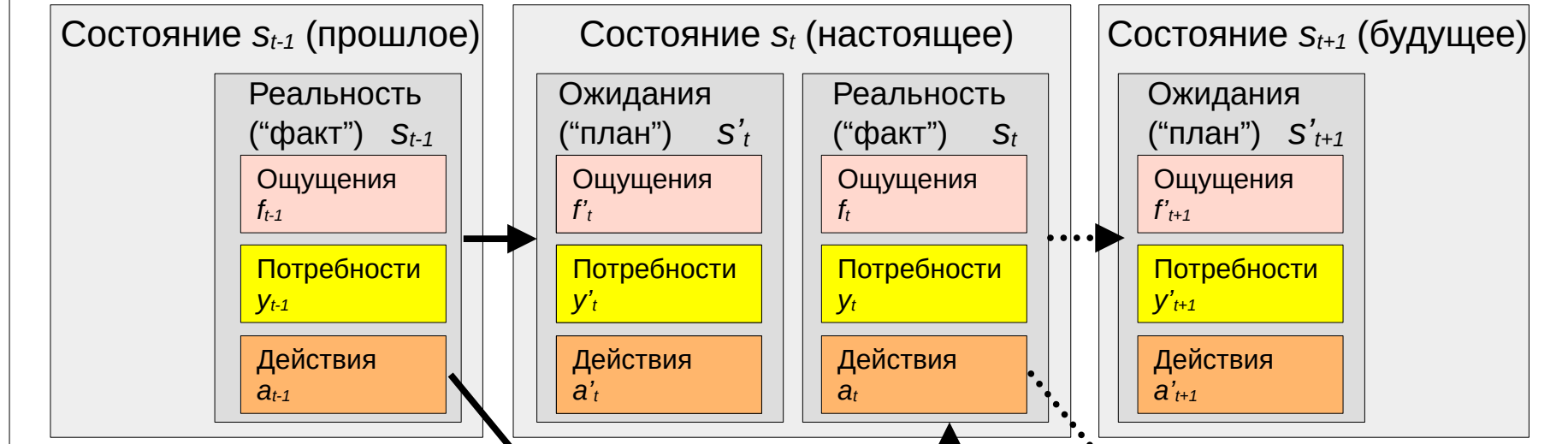
Интеллект = Система принятия решений

Модели  $s$  (“инварианты”) состояний, полезности  $U$  и вероятности  $P$  переходов  
 $U(\{s_t\}_{t \in \{-T, -1\}}, s'_0) = L(x \cdot (y_t - y_{t+1}), s'_t - s_t, E(a_t)) \quad s'_t = \operatorname{argmax}_s (U(\{s_t\}_{t \in \{-T, -1\}}, s'_t), P(\{s_t\}_{t \in \{-T, -1\}}, s'_t))$

↑ Обучение на опыте

↓ Принятие решений

Пространство состояний и эпизодическая память (“прецеденты”)



А.Г.Колонин, В.Г.Крюков:  
Вычислительная концепция психики,  
Статья принята на конференцию  
Нейроинформатика-25

Решения

Актуаторы

Ощущения

Сенсоры

DQN?

Трансформеры на токенизированных  
состояниях?

“Глобальное подкрепление”!



# Психика = Операционная система

## Интеллект = Система принятия решений

Модели  $s$  ("инварианты") состояний, полезности  $U$  и вероятности  $P$  переходов

$$U(\{s_t\}_{t \in \{-T, -1\}}, s'_0) = L(x \cdot (y_t - y_{t+1}), s'_t - s_t, E(a_t)) \quad s'_t = \operatorname{argmax}_s (U(\{s_t\}_{t \in \{-T, -1\}}, s'_t), P(\{s_t\}_{t \in \{-T, -1\}}, s'_t))$$

↑ Обучение на опыте

↓ Принятие решений

Пространство состояний и эпизодическая память ("прецеденты")

Состояние  $s_{t-1}$  (прошое)

Реальность  
("факт")  $s_{t-1}$

Ощущения  
 $f_{t-1}$

Потребности  
 $y_{t-1}$

Действия  
 $a_{t-1}$

Состояние  $s_t$  (настоящее)

Ожидания  
("план")  $s'_t$

Ощущения  
 $f'_t$

Потребности  
 $y'_t$

Действия  
 $a'_t$

Реальность  
("факт")  $s_t$

Ощущения  
 $f_t$

Потребности  
 $y_t$

Действия  
 $a_t$

Состояние  $s_{t+1}$  (будущее)

Ожидания  
("план")  $s'_{t+1}$

Ощущения  
 $f'_{t+1}$

Потребности  
 $y'_{t+1}$

Действия  
 $a'_{t+1}$

Решения

Ощущения

Актуаторы

Сенсоры

$x \cdot y_t$  — "мотивационный вектор"  
V. F. Petrenko and A. P. Suprun, "Goal oriented systems, evolution, and the subjective aspect in systemology," Tr. Inst. Sistem. Analiza RAN 62 (1) (2012)

Copyright © 2026 Anton Kolonin, Aigents®

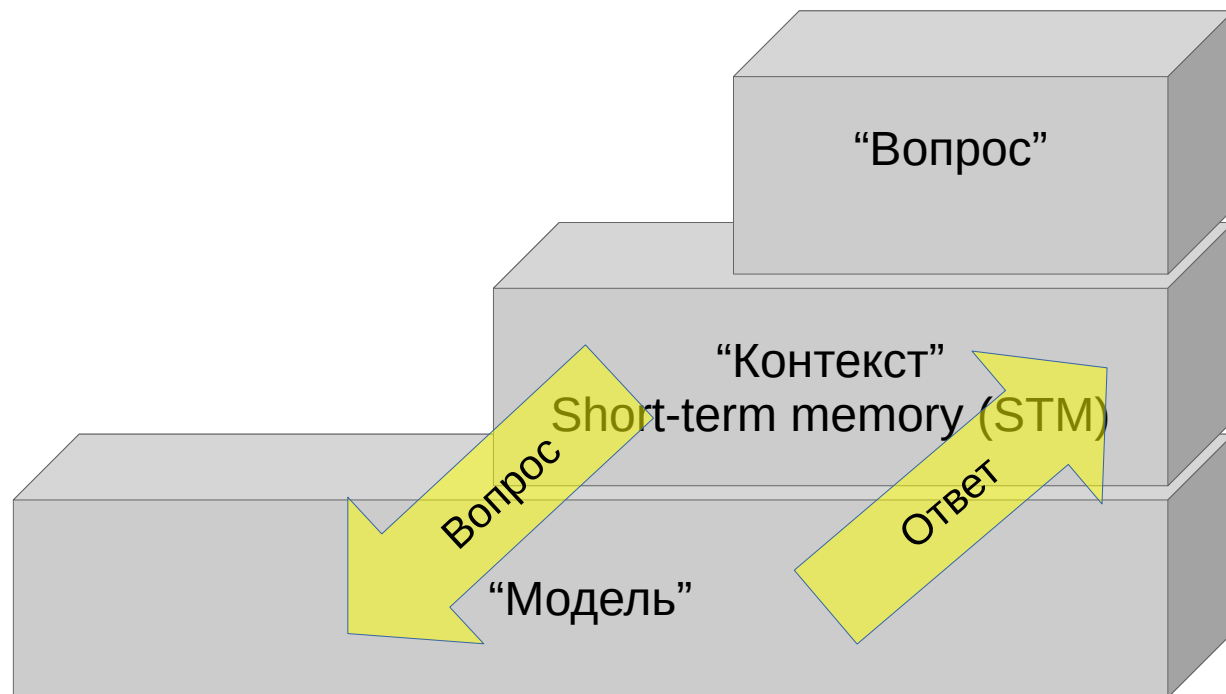
Кластеризация состояний?  
Асинхронная их обработка?

# Обучение модели (БЯМ) - Training/Learning



*Модель замораживается после обучения*

# Использование модели (БЯМ) - Inference

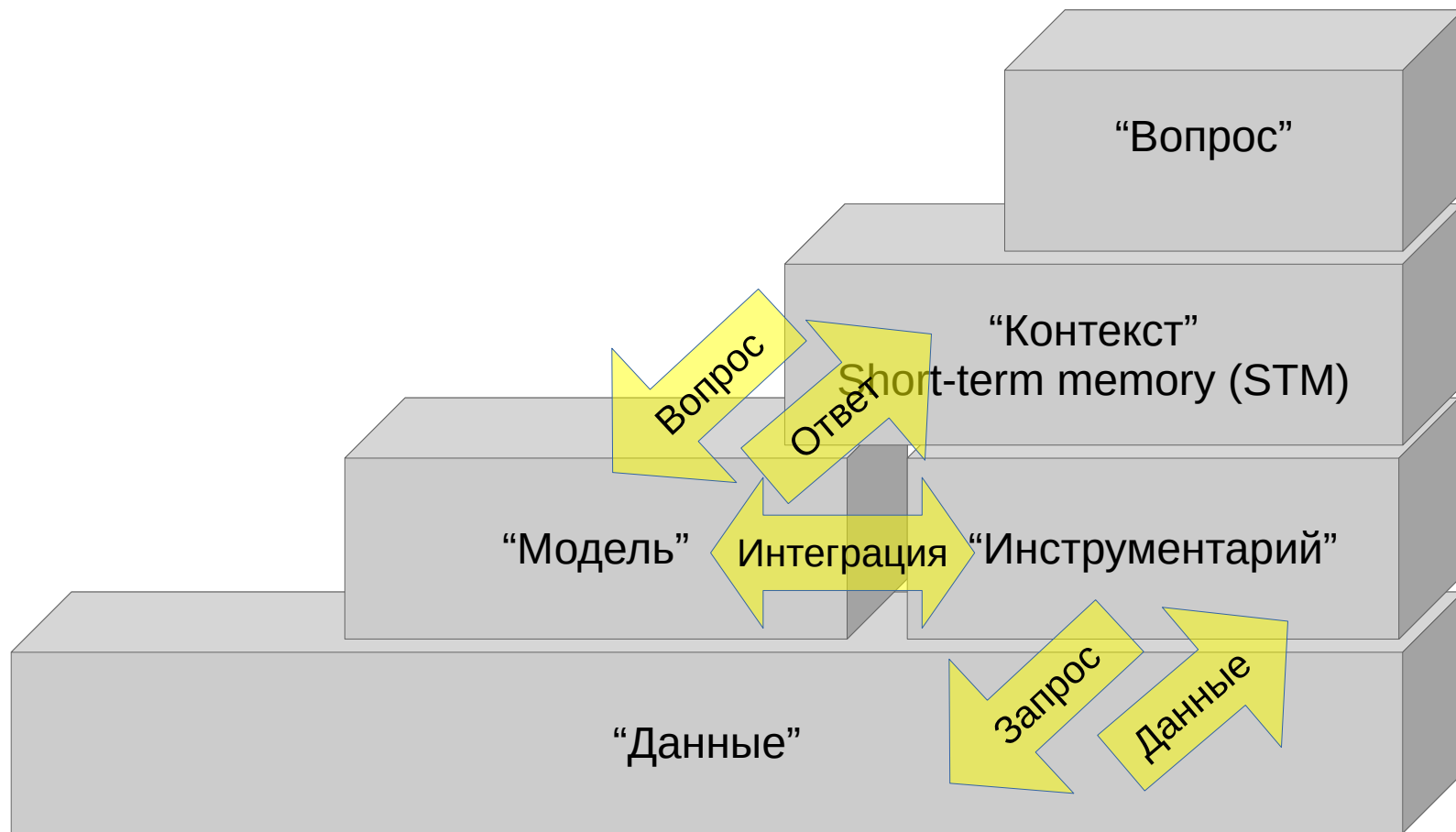


*Данные не участвуют в выводе*

*Новые данные и контексты требуют переобучения на старых и новых*

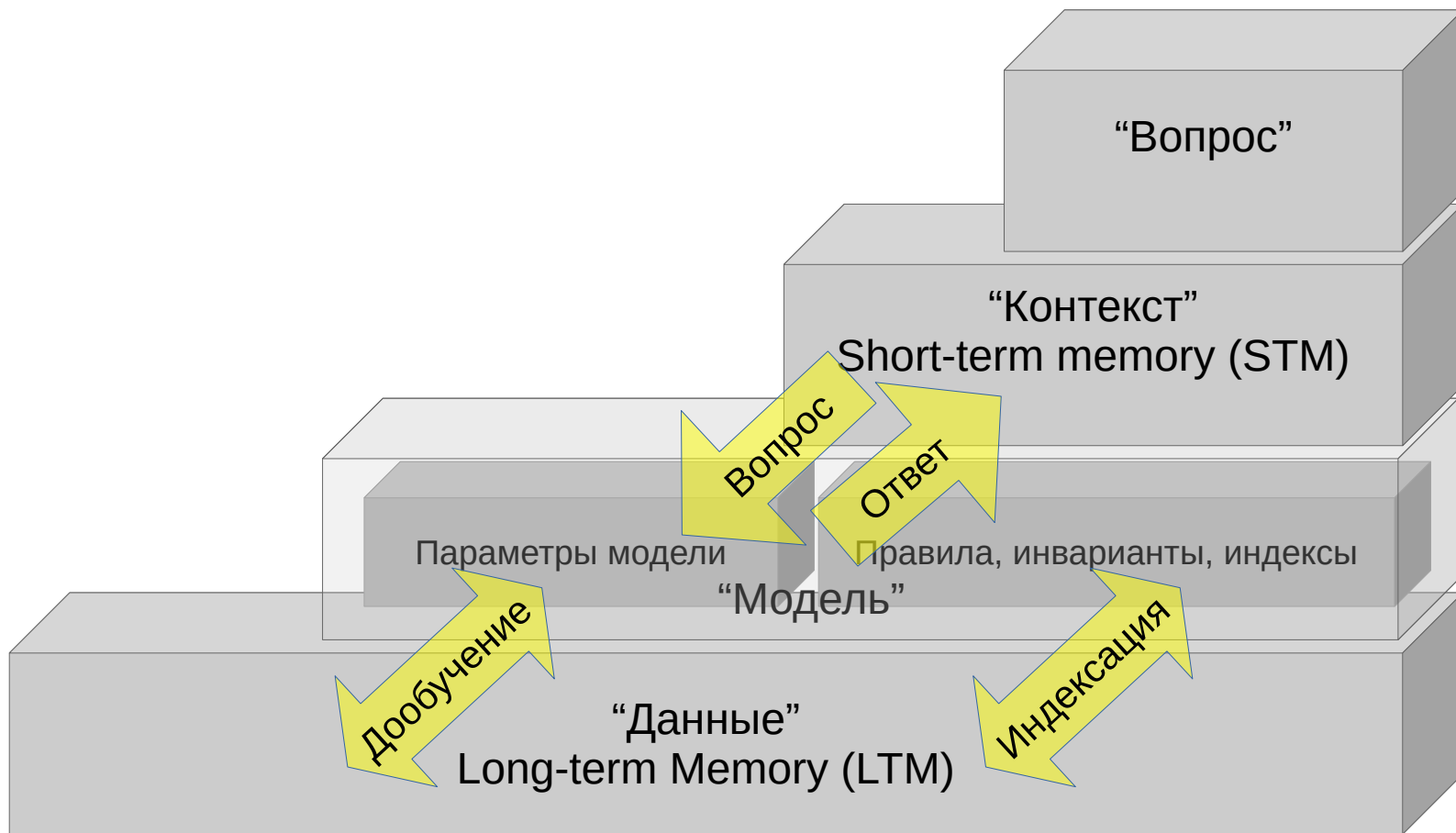
*На это все может не хватить параметров конкретной модели*

# Много-агентные модели (БЯМ+RAG/GraphRAG)



*Данные могут подкачиваться в контекст и обновляться с его учетом  
Модель все равно требует переобучения с учетом новых данных*

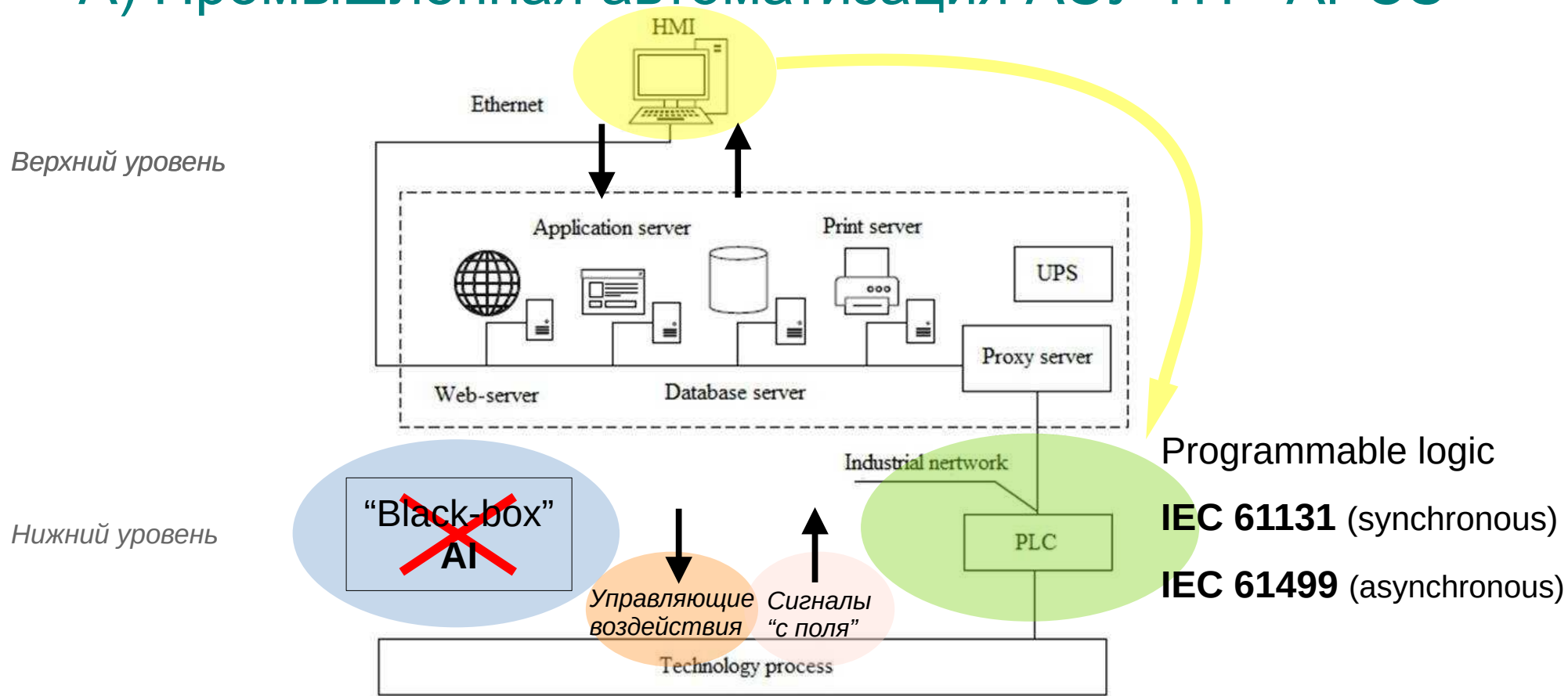
# Гибридная архитектура с долгосрочной памятью



*Возможно инкрементальное "дообучение" на динамически пополняющихся данных  
Динамически обновляющиеся данные могут инкрементально обновлять соответствующие "инварианты"*

# Применение – Use case:

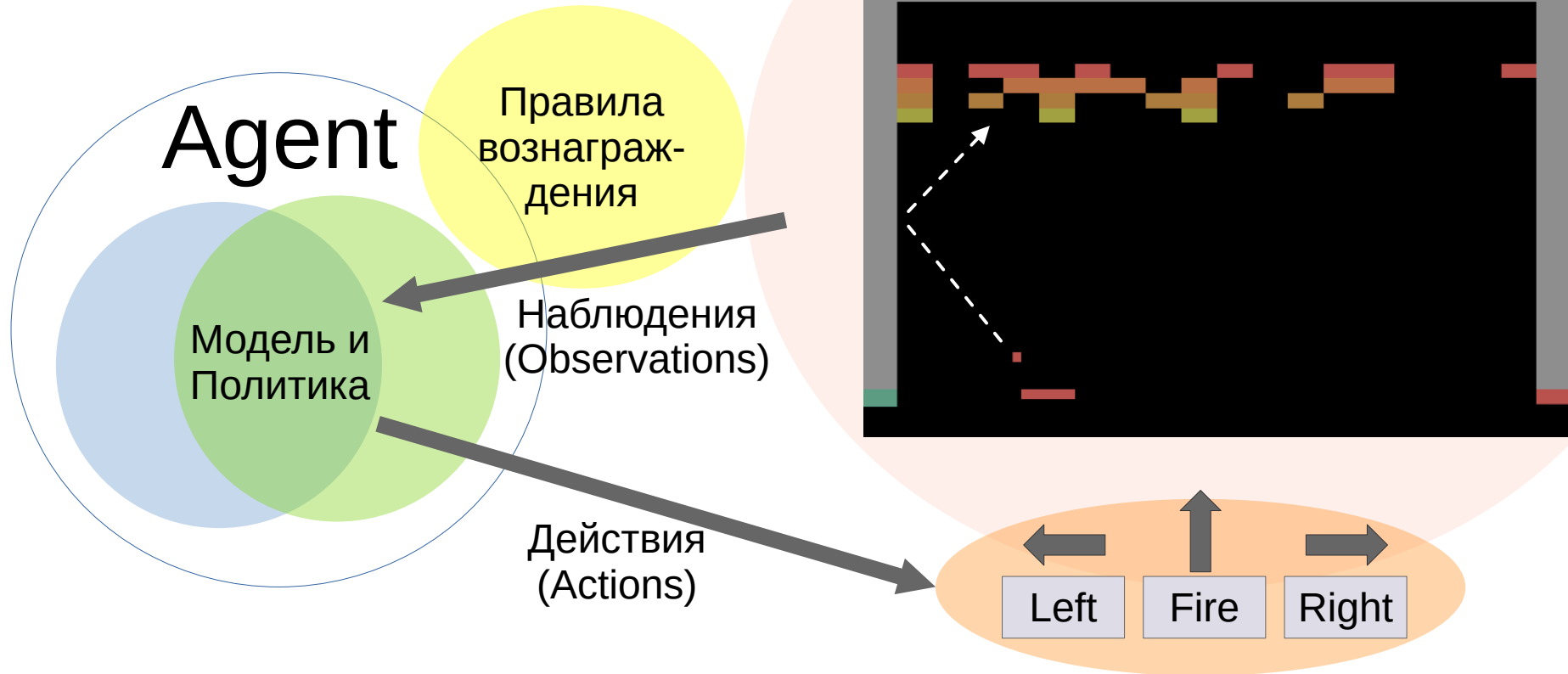
## А) Промышленная автоматизация АСУ ТП - APCS



[https://www.researchgate.net/publication/311662442\\_Adaptive\\_Intelligent\\_Manufacturing\\_Control\\_Systems](https://www.researchgate.net/publication/311662442_Adaptive_Intelligent_Manufacturing_Control_Systems)

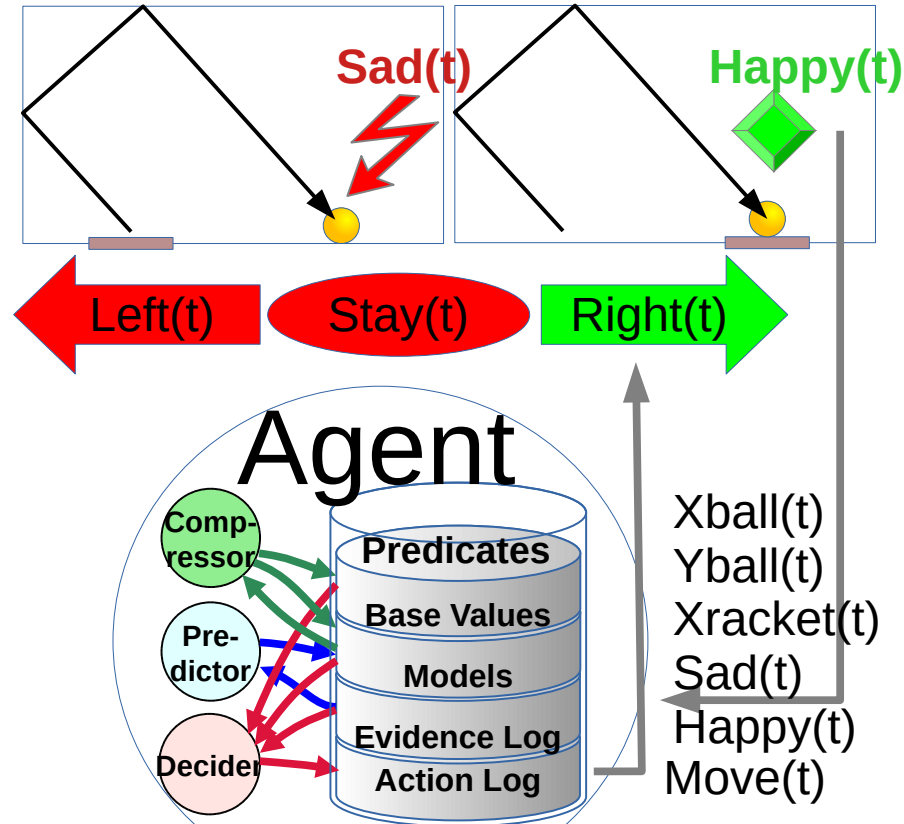
# Применение – Use case:

## В) Виртуальная игровая среда OpenAI Gym (Atari Breakout)

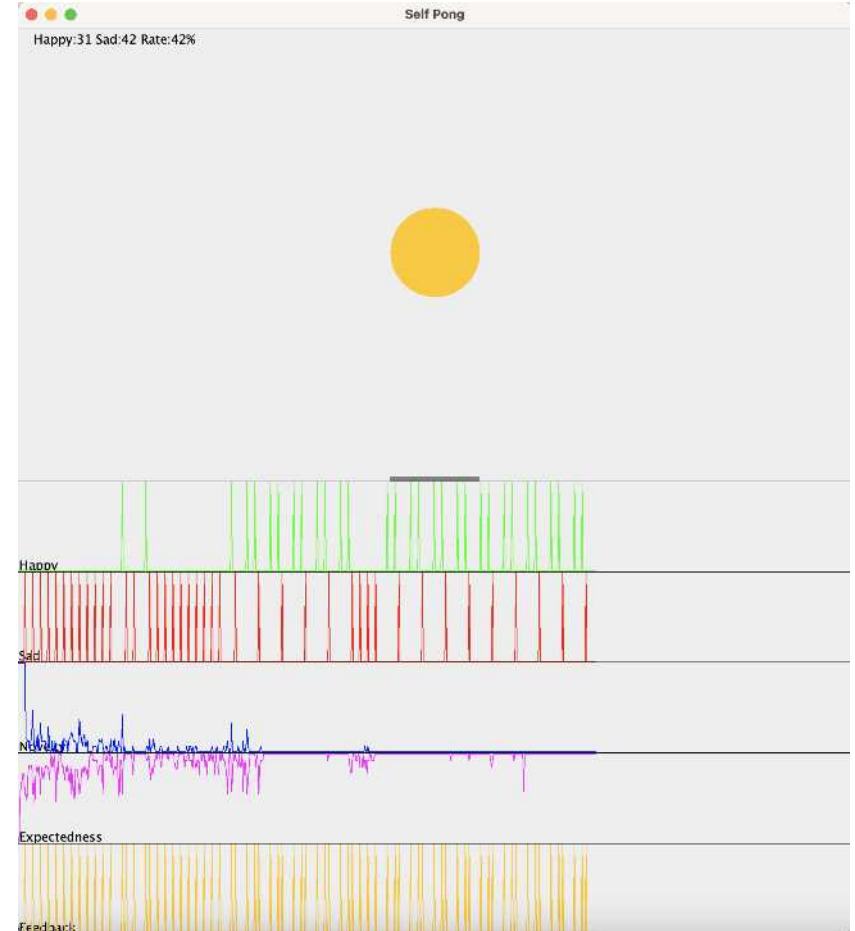


# Когнитивная архитектура на графах последовательных состояний

## State-based History-aware Artificial Reinforcement Intelligence Kernel (**Sharik**)

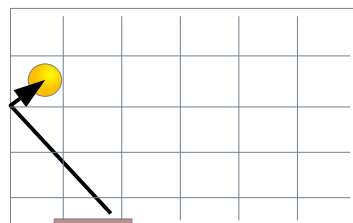
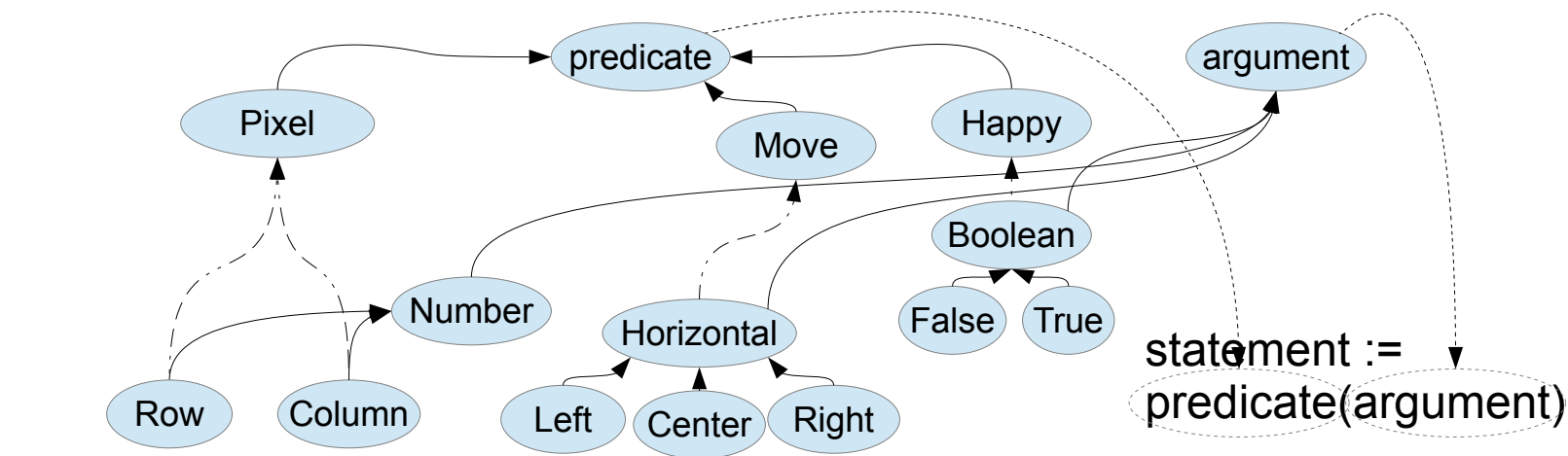


А.Г.Колонин, В.Г.Крюков:  
Вычислительная концепция психики,  
Статья подана на конференцию  
Нейроинформатика-25

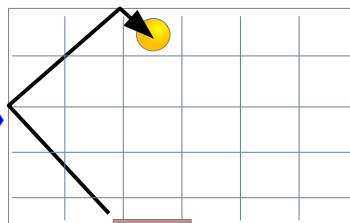




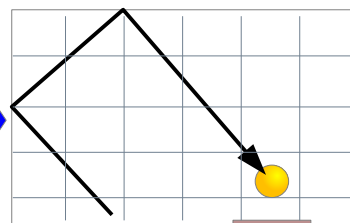
# Играем в “1 player Ping-Pong” – уровень пикселей



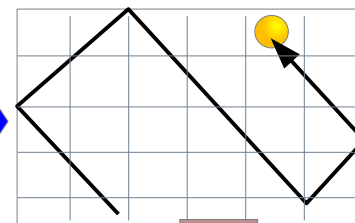
Pixel(1,0)  
Pixel(4,1)  
Happy(False)  
=>  
**Move(Left)**



Pixel(0,2)  
Pixel(4,2)  
Happy(False)  
=>  
**Move(Right)**



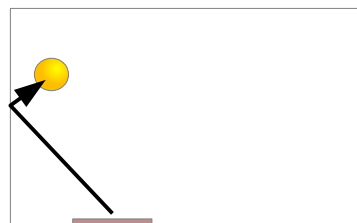
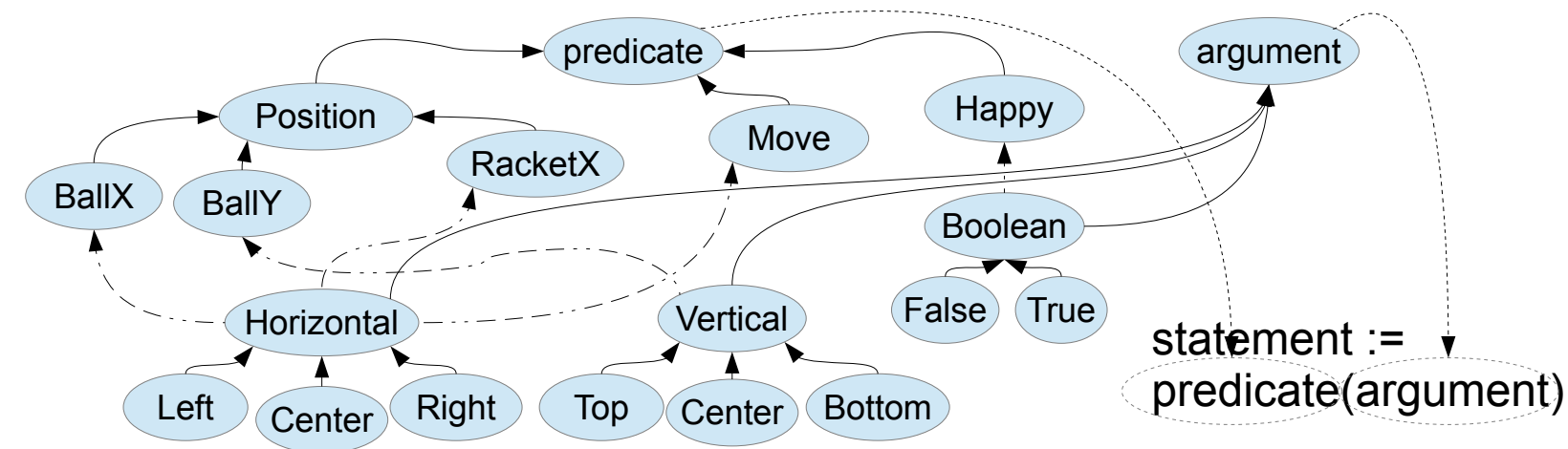
Pixel(3,4)  
Pixel(4,4)  
Happy(False)  
=>  
**Move(Left)**



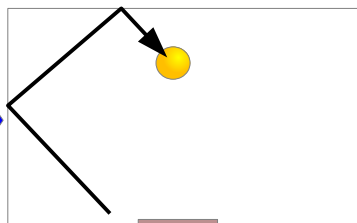
Pixel(0,4)  
Pixel(4,3)  
**Happy(True)**  
=>  
Move(Left)

Global  
Feedback

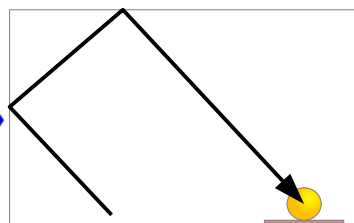
# Играем в “1 player Ping-Pong” – уровень объектов



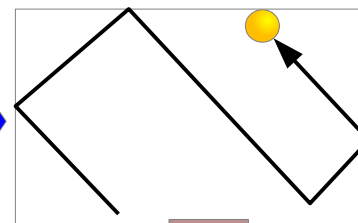
BallY(Top)  
BallX(Left)  
RacketX(Left)  
Happy(False)  
=> **Move(Left)**



BallY(Top)  
BallX(Center)  
RacketX(Center)  
Happy(False)  
=> **Move(Right)**



BallY(Bottom)  
BallX(Right)  
RacketX(Right)  
Happy(False)  
=> **Move(Right)**

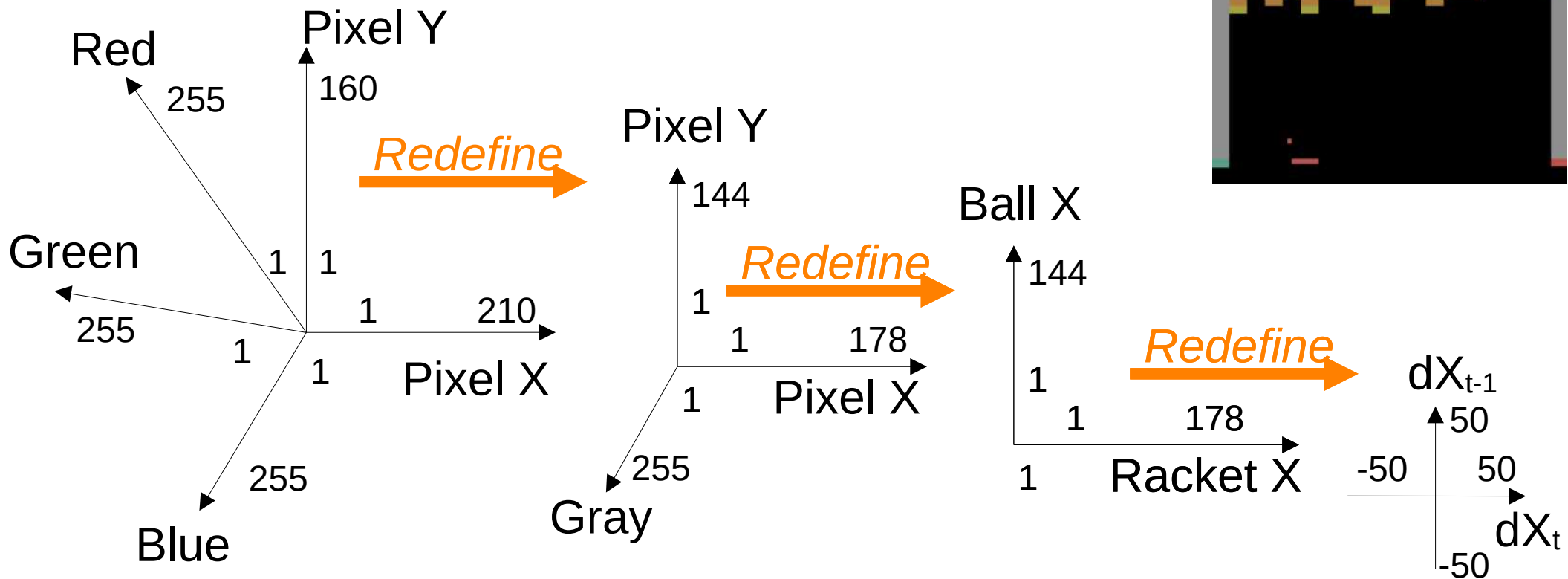


BallY(Bottom)  
BallX(Right)  
RacketX(Right)  
**Happy(True)**  
=> **Move(Left)**

Global  
Feedback

# Проблемы (понижения) размерности и (повышения) дискретности

Re-defining environment in Atari Breakout



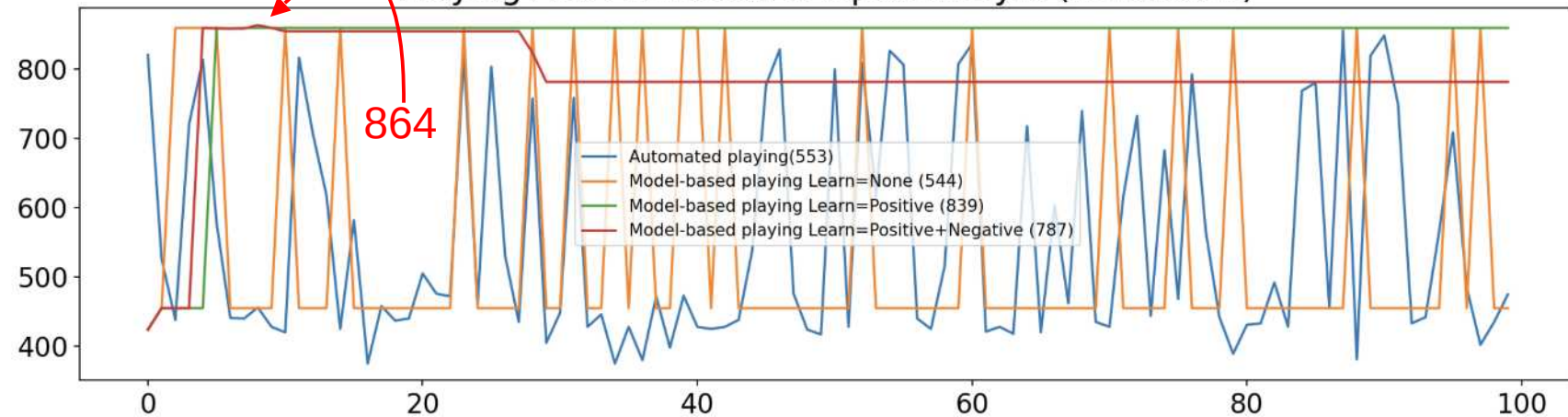
Interpretable representation learning for 3D multi-piece intracellular structures using point clouds

<https://www.nature.com/articles/s41592-025-02729-9>

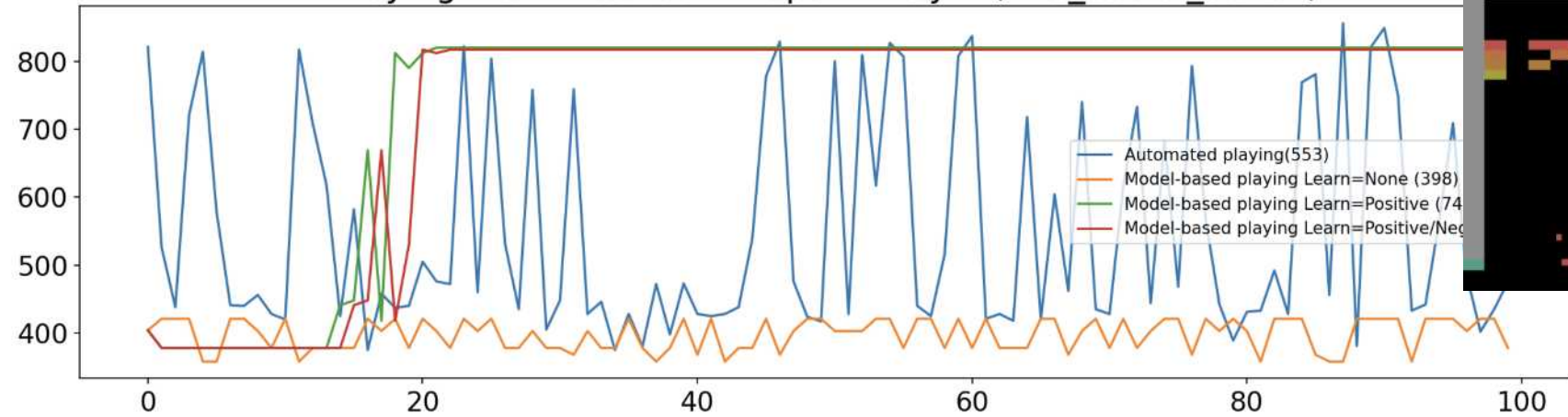
Copyright © 2026 Anton Kolonin, Aigents®

# “Imitation learning” – принятие решений на “пред-обученной” модели

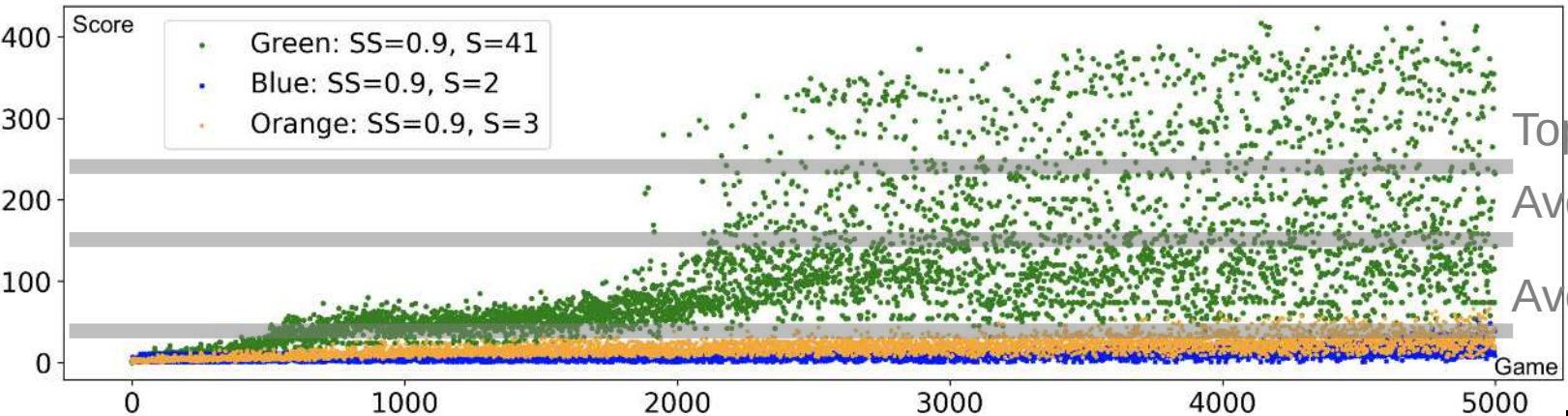
Playing Atari Breakout in Open AI Gym (Nov32025)



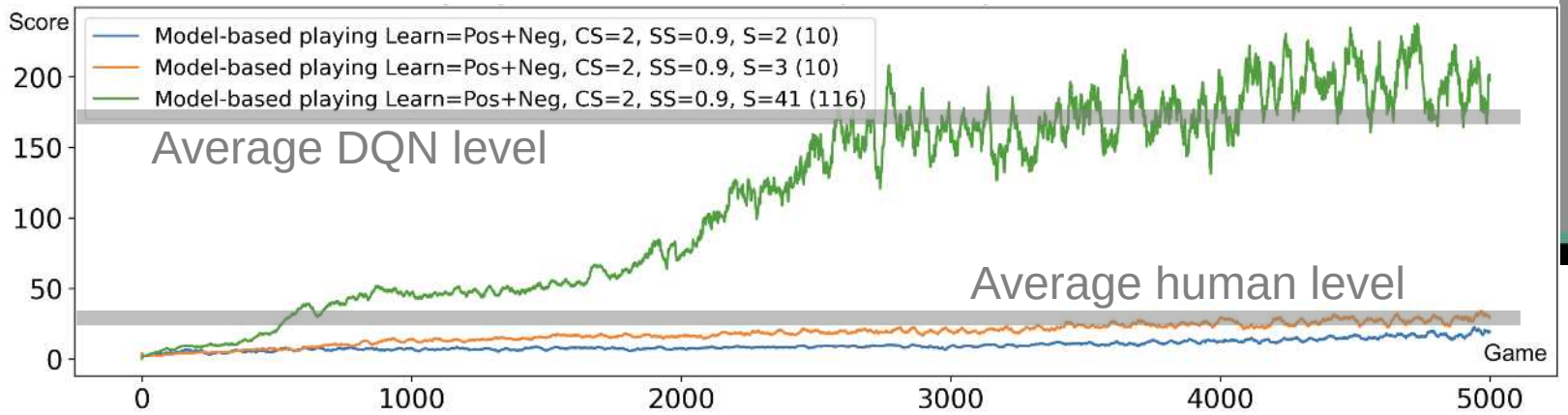
Playing Atari Breakout in Open AI Gym (find\_useful\_action)



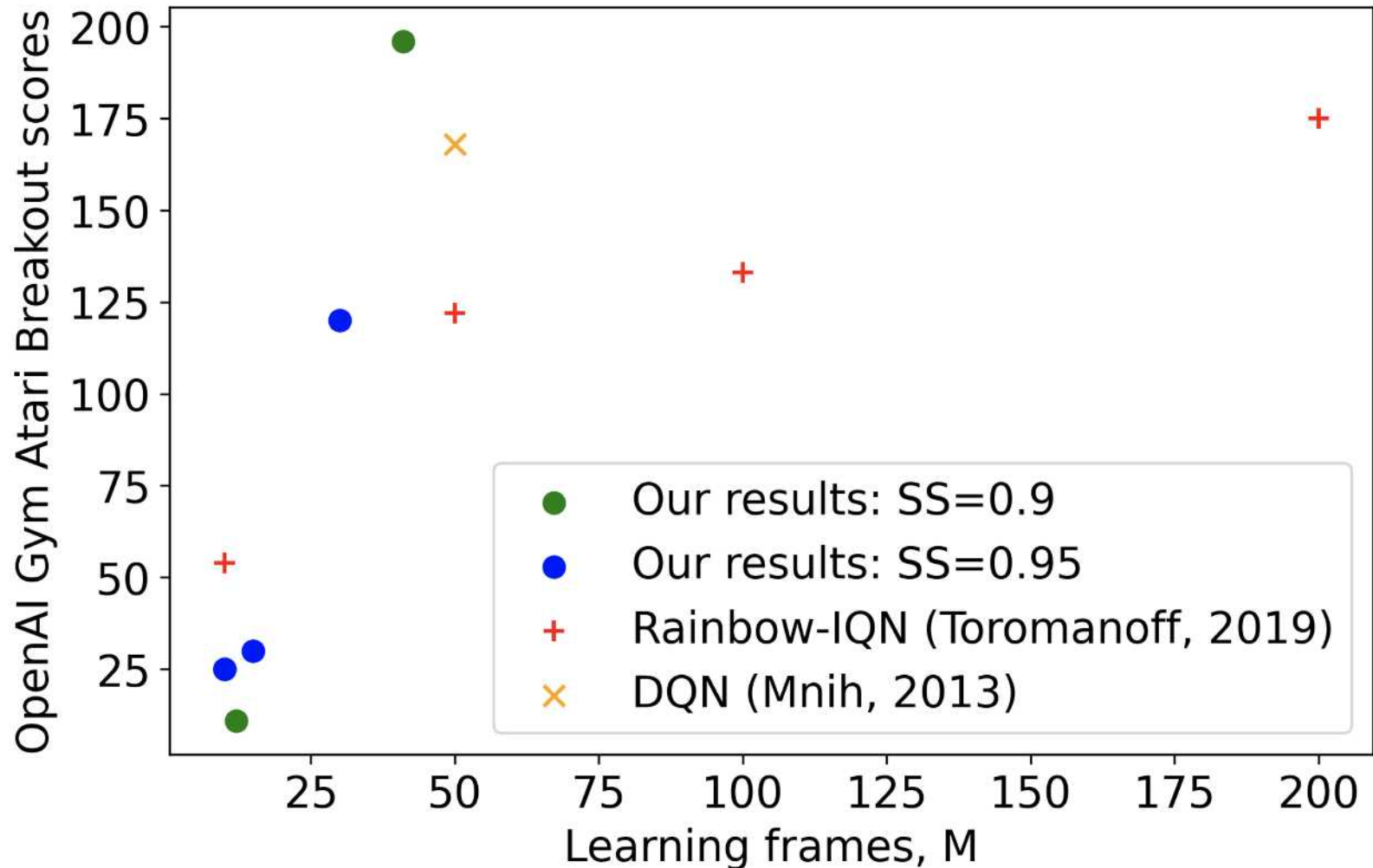
# Обучение на собственном опыте – переходы между последовательностями



Top DQN level  
Average DQN level  
Average human level



# Обучение на собственном опыте – сравнение с “baselines”

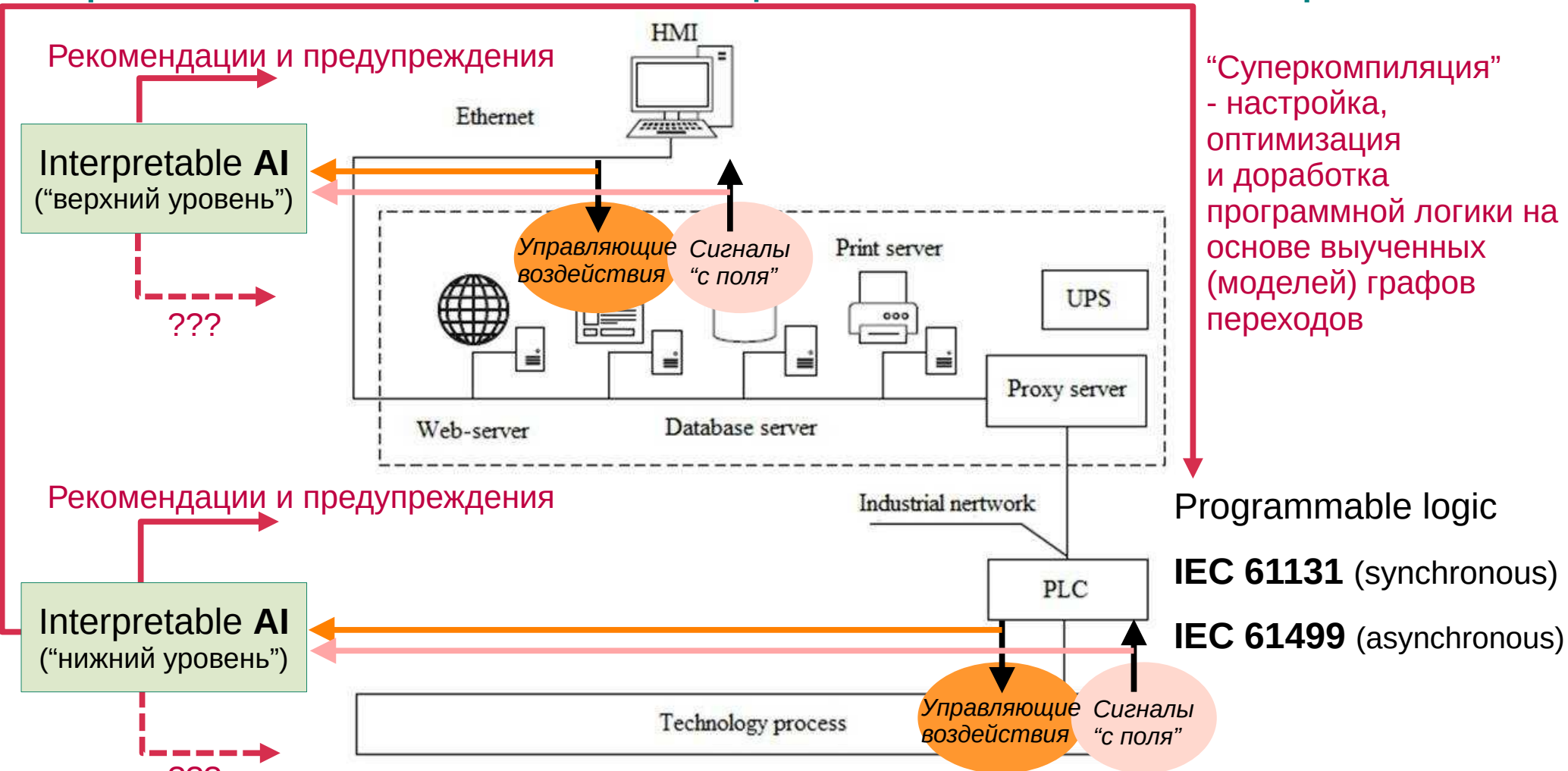


# Что дальше?

1. Стабилизация обучаемости!
2. Интерпретируемое понижение размерности!
3. Больше окружений – хороших и разных!
4. Кластеризация/сегментация пространства состояний и  
раздельное исполнение сегментов (*пинг-понг по-македонски*)?
5. Прикладное применение – промышленная автоматизация?
6. Поддержка иерархий/гетерархий пространств состояний?\*
7. Формализация перевода графов состояний в языки  
программной логики (“суперкомпиляция”)?\*\*

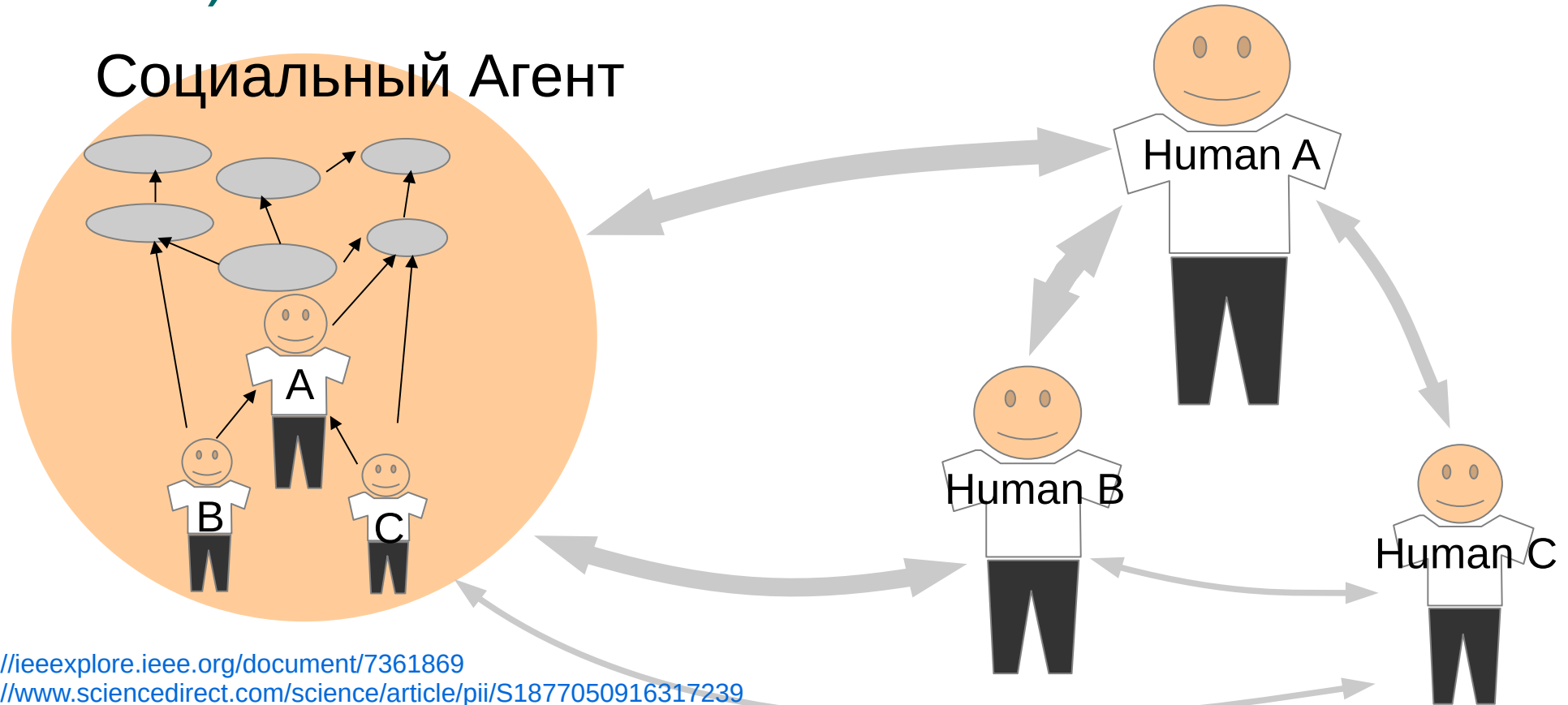


# Промышленная автоматизация – постановка эксперимента





# Система поддержки принятия решений (СППР) на основе коллективного интеллекта

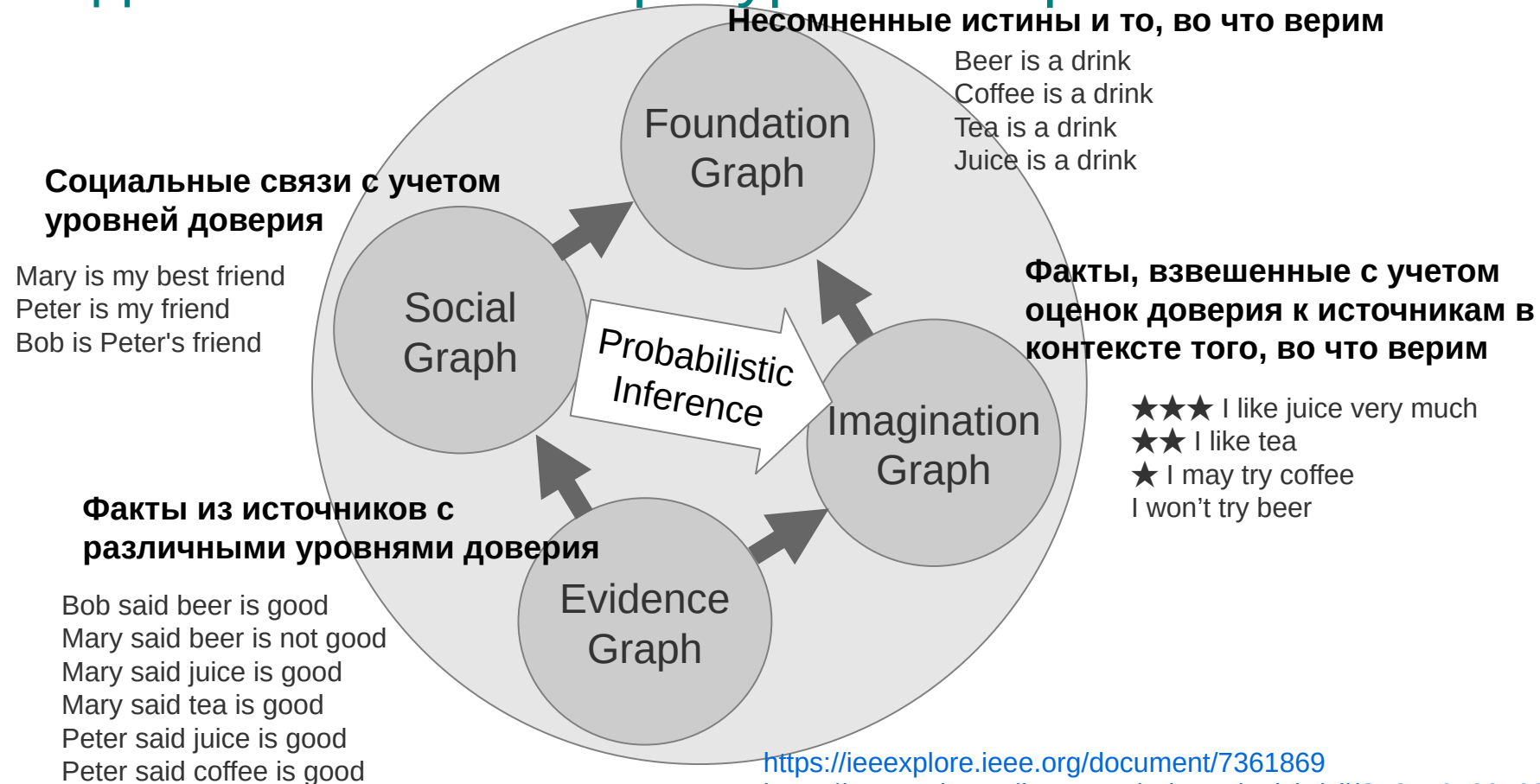


<https://ieeexplore.ieee.org/document/7361869>

<https://www.sciencedirect.com/science/article/pii/S1877050916317239>

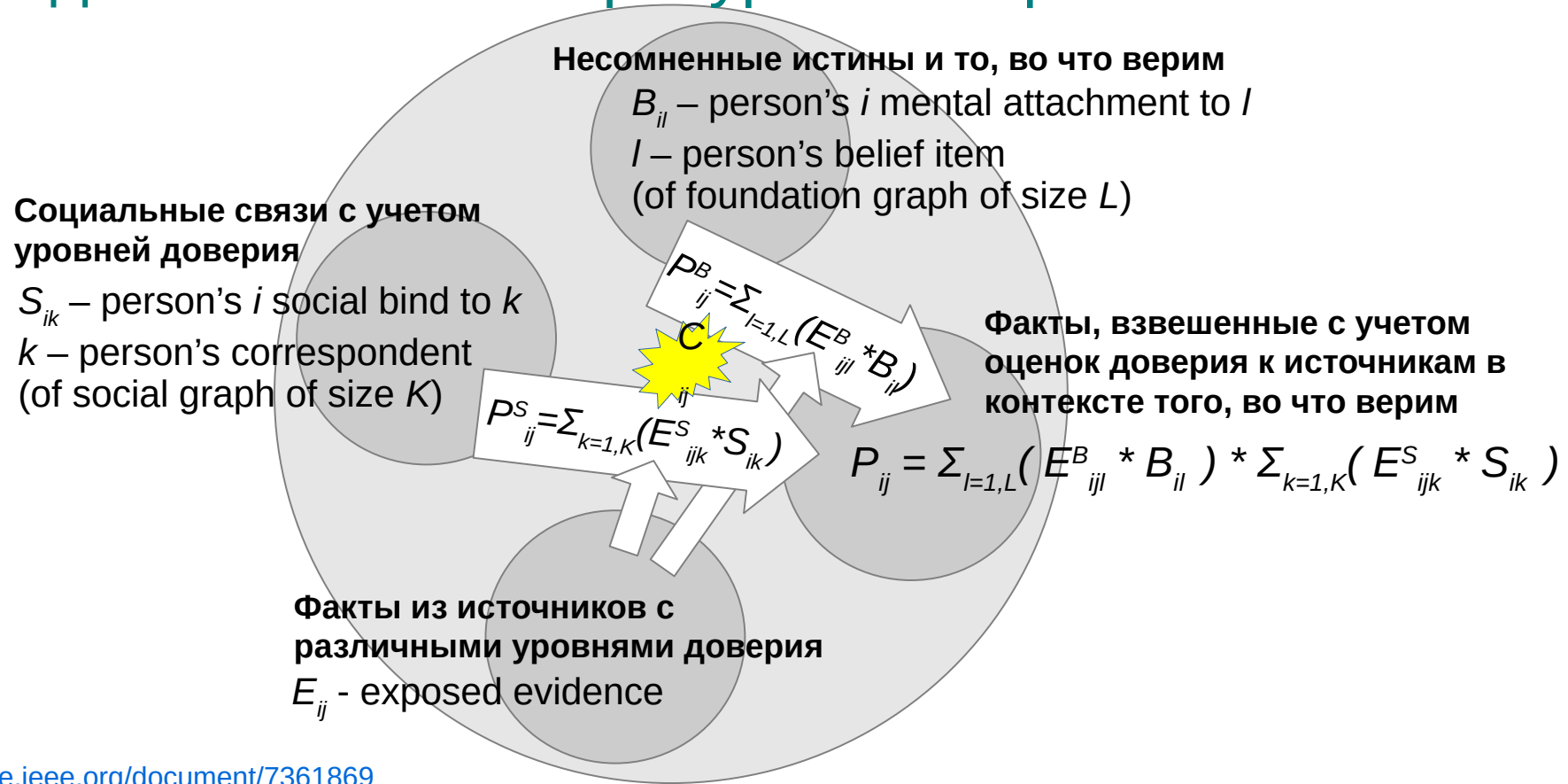
[https://link.springer.com/chapter/10.1007/978-3-319-97676-1\\_10](https://link.springer.com/chapter/10.1007/978-3-319-97676-1_10)

# Когнитивно-поведенческая модель на основе социального доказательства и ресурсного ограничения



<https://ieeexplore.ieee.org/document/7361869>  
<https://www.sciencedirect.com/science/article/pii/S1877050916317239>  
[https://link.springer.com/chapter/10.1007/978-3-319-97676-1\\_10](https://link.springer.com/chapter/10.1007/978-3-319-97676-1_10)

# Когнитивно-поведенческая модель на основе социального доказательства и ресурсного ограничения



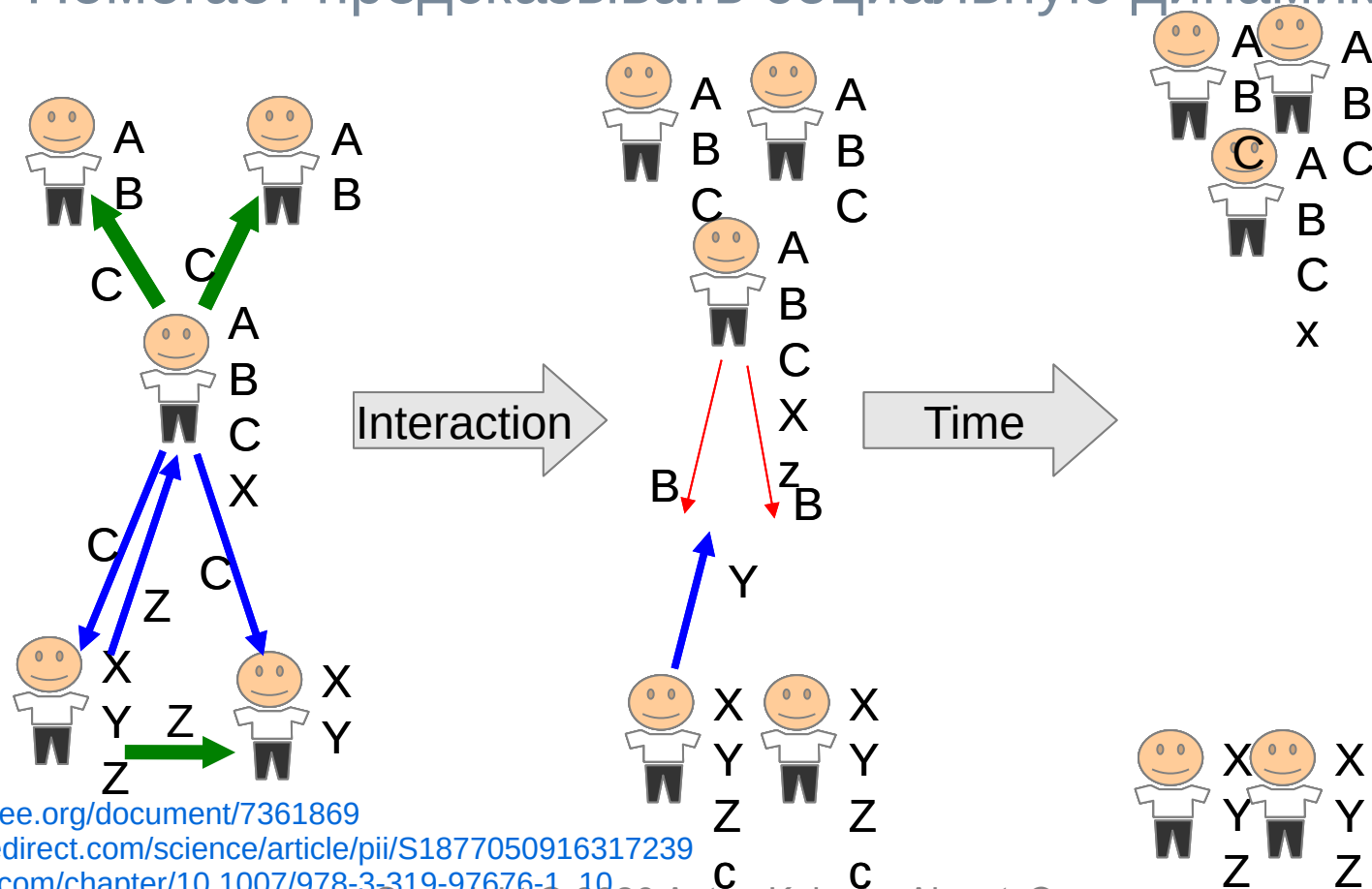
<https://ieeexplore.ieee.org/document/7361869>

<https://www.sciencedirect.com/science/article/pii/S1877050916317239>

[https://link.springer.com/chapter/10.1007/978-3-319-97676-1\\_10](https://link.springer.com/chapter/10.1007/978-3-319-97676-1_10)

# Когнитивно-поведенческая модель на основе социального доказательства и ресурсного ограничения

Помогает предсказывать социальную динамику



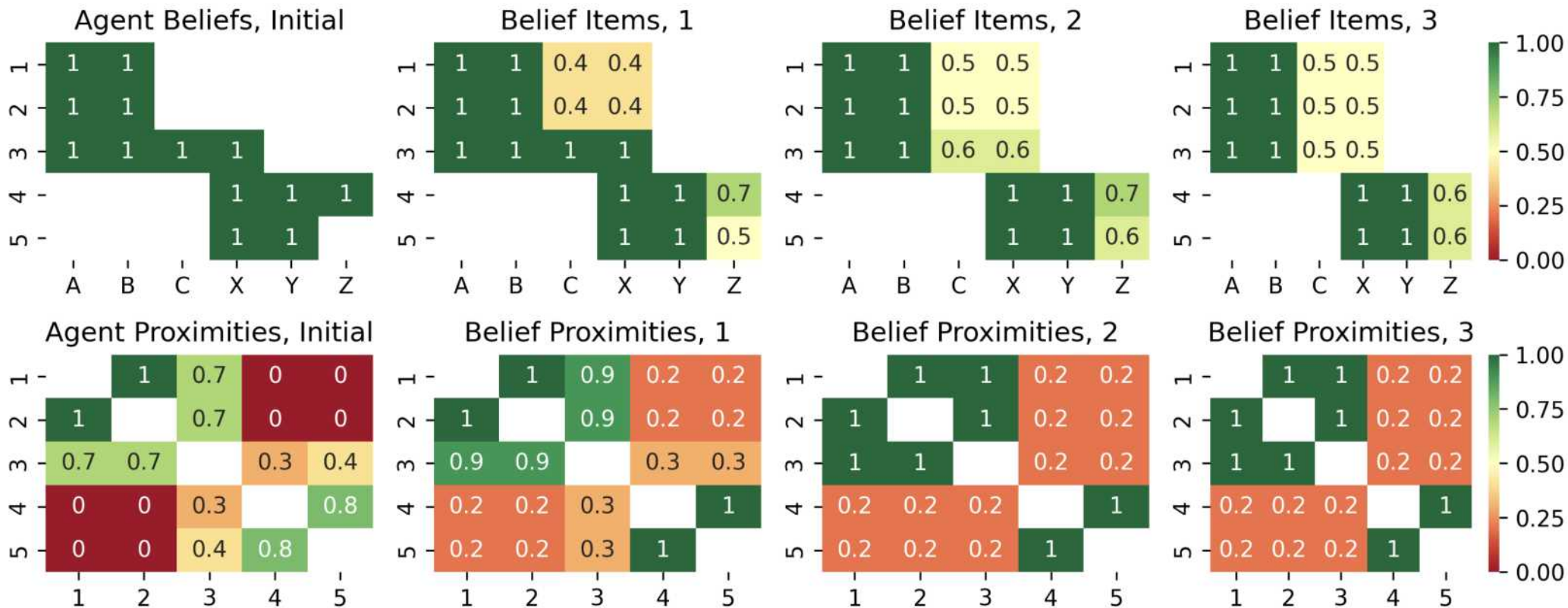
<https://ieeexplore.ieee.org/document/7361869>

<https://www.sciencedirect.com/science/article/pii/S1877050916317239>

[https://link.springer.com/chapter/10.1007/978-3-319-97676-1\\_10](https://link.springer.com/chapter/10.1007/978-3-319-97676-1_10)

Copyright © 2026 Anton Kolonin, Aigents®

# Симуляция: социальная ограниченность, когнитивный ресурс



Results of three rounds of multi-agent simulation with the forgetting threshold peer threshold = 0.5, the social relationship threshold forgetting threshold = 0.0 (limiting social connections K without constraint on belief capacity L). Top row: values of belief items A, B, C, X, Y, Z for five agents 1, 2, 3, 4, 5. Bottom row: social proximity matrices between agents, estimated as the cosine similarity of their beliefs. From left to right: the initial state before the simulation, and then the updated states of beliefs and proximities after three subsequent rounds of multi-agent communication.

# Спасибо за внимание! Вопросы?

Антон Колонин  
[akolonin@aigents.com](mailto:akolonin@aigents.com)  
Telegram: [akolonin](#)

Запись семинара по  
теме доклада



Статья по теме доклада,  
принятая на конференцию  
Нейроинформатика-2025

