# Distributed Algorithms as a Gateway to Thinking Slow

*Roger Wattenhofer*

# Deep Learning is Robust to Massive Label Noise

David Rolnick [*1]   Andreas Veit [*2]   Serge Belongie [2]   Nir Shavit [3]

# Machine Learning with Adversaries:
# Byzantine Tolerant Gradient Descent

**Peva Blanchard**
EPFL, Switzerland
peva.blanchard@epfl.ch

**El Mahdi El Mhamdi**[*]
EPFL, Switzerland
elmahdi.elmhamdi@epfl.ch

**Rachid Guerraoui**
EPFL, Switzerland
rachid.guerraoui@epfl.ch

**Julien Stainer**
EPFL, Switzerland
julien.stainer@epfl.ch

# QSGD: Communication-Efficient SGD
# via Gradient Quantization and Encoding

**Dan Alistarh**
IST Austria & ETH Zurich
dan.alistarh@ist.ac.at

**Demjan Grubic**
ETH Zurich & Google
demjangrubic@gmail.com

**Jerry Z. Li**
MIT
jerryzli@mit.edu

**Ryota Tomioka**
Microsoft Research
ryoto@microsoft.com

**Milan Vojnovic**
London School of Economics
M.Vojnovic@lse.ac.uk

# Byzantine Fault-Tolerant Distributed Machine Learning using D-SGD and Norm-Based Comparative Gradient Elimination (CGE)

Nirupam Gupta
*EPFL*
Lausanne, Switzerland
nirupam115@gmail.com

Shuo Liu
*Georgetown University*
Washington, D.C., USA
sl1539@georgetown.edu

Nitin Vaidya
*Georgetown University*
Washington, D.C., USA
nv198@georgetown.edu

# Distributed Algorithms as a Gateway to Thinking Slow
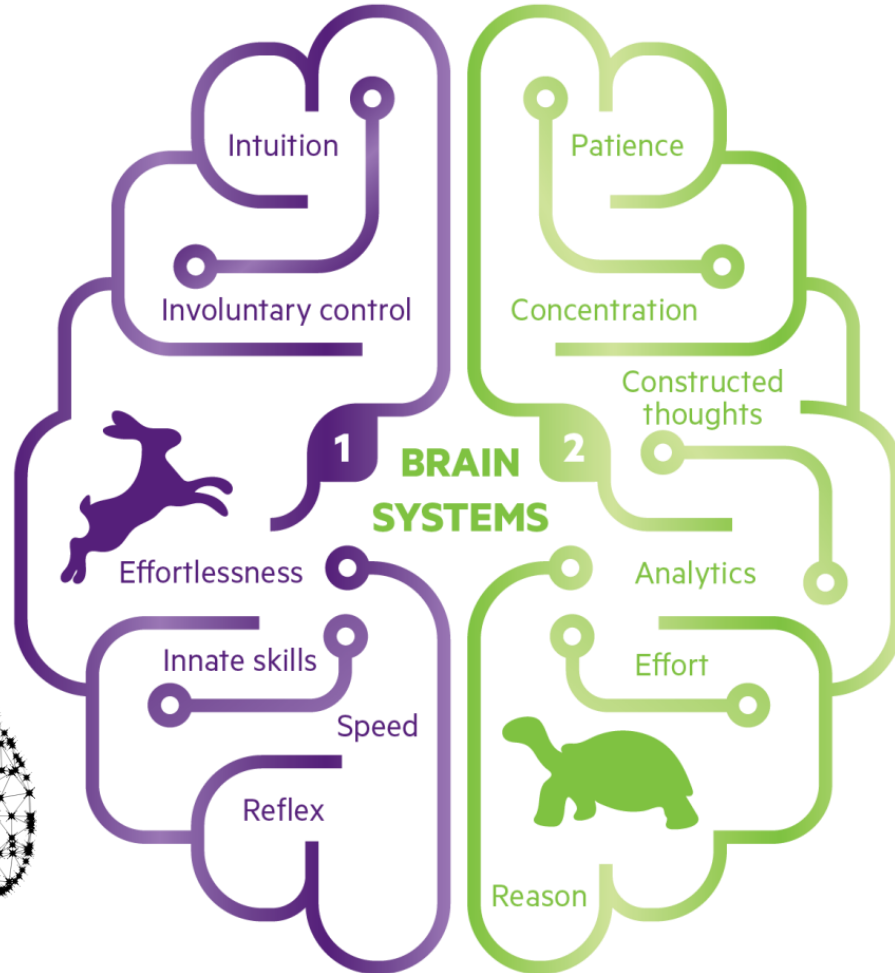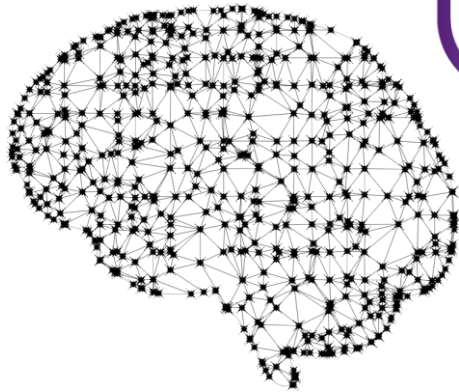
*Roger Wattenhofer*

# THINKING,
# FAST and SLOW

# DANIEL
# KAHNEMAN

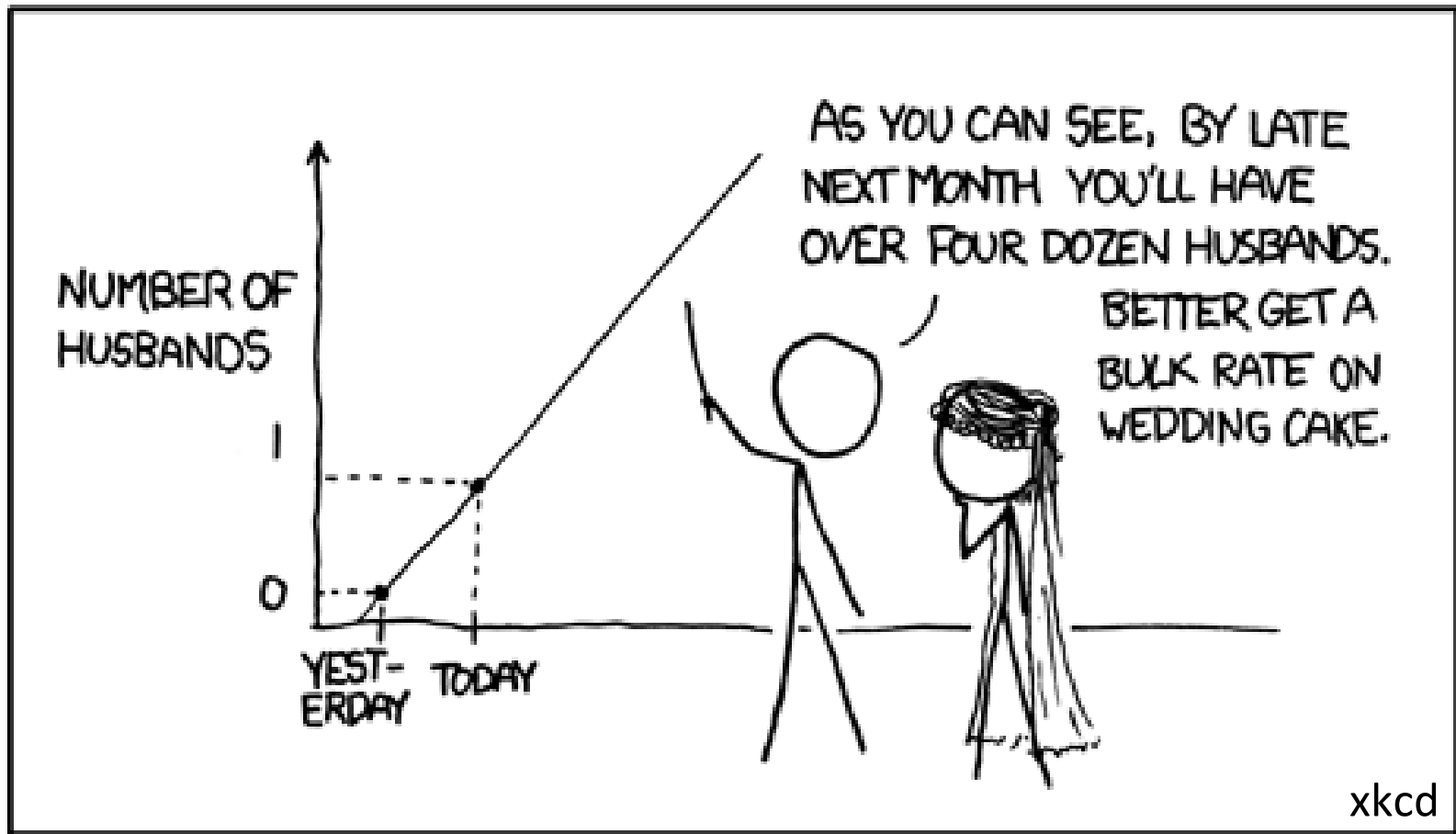WINNER OF THE NOBEL PRIZE IN ECONOMICS

# Classification

# Chihuahua
# vs.
# Muffin
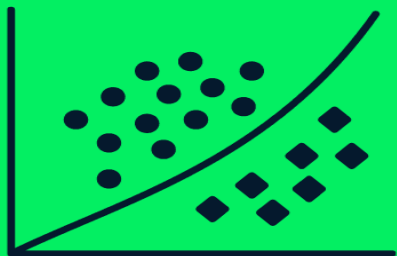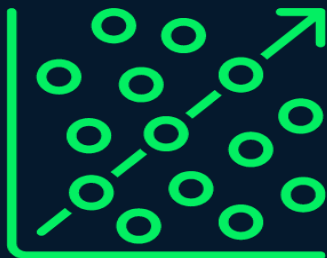
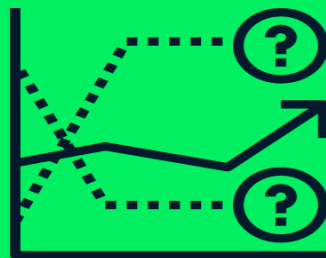Linear Regression

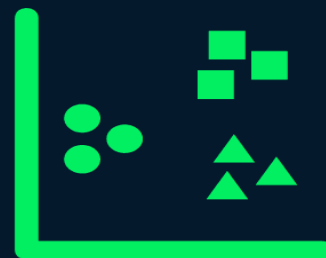Classification

Regression

**Classification**

**Regression**

**Forecast**
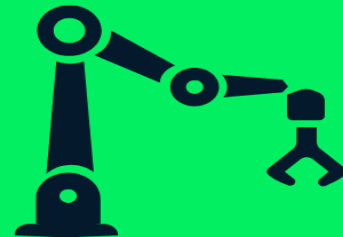
**Clustering**

**Machine Translation**

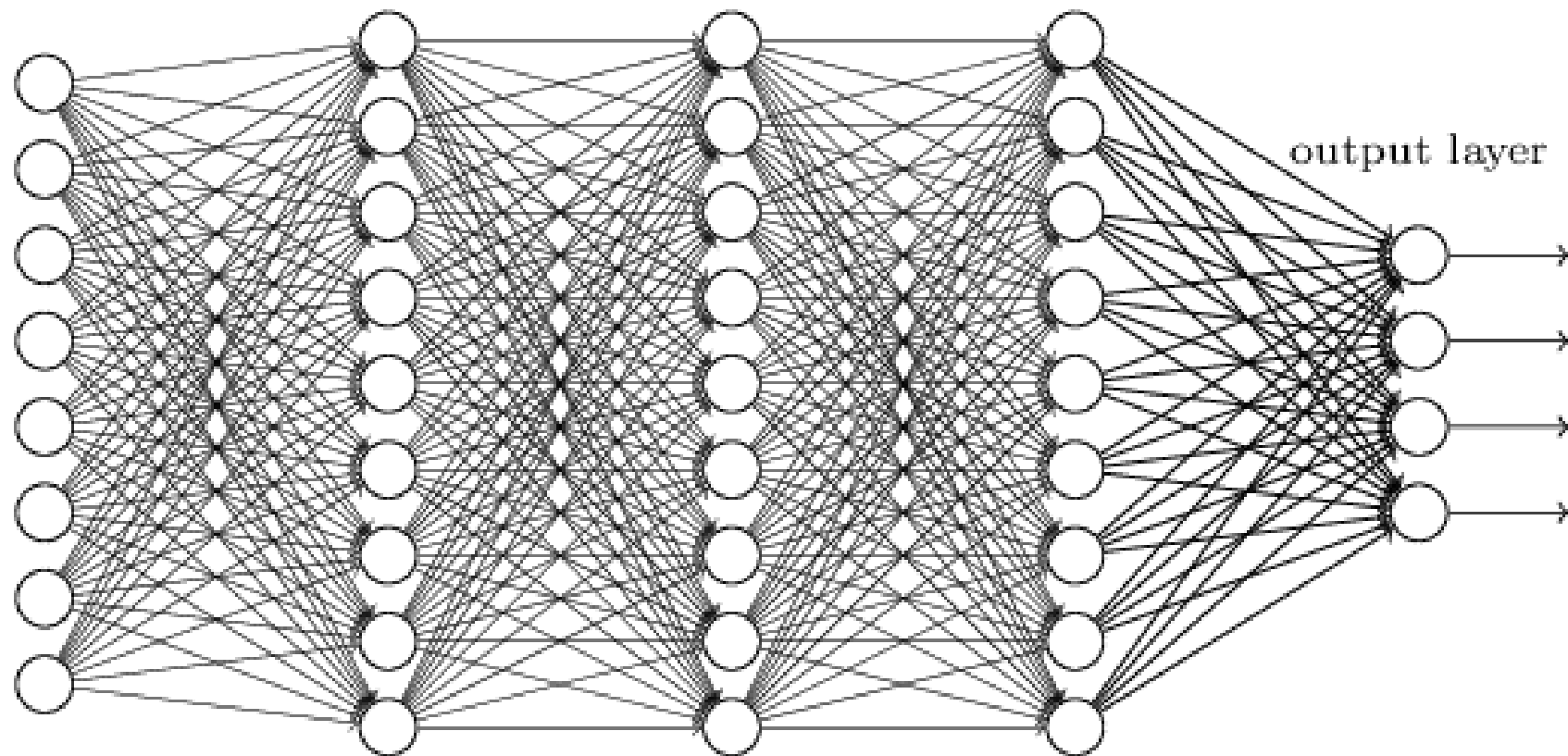**Computer Vision**

**Generative Art**

**Reinforcement learning**

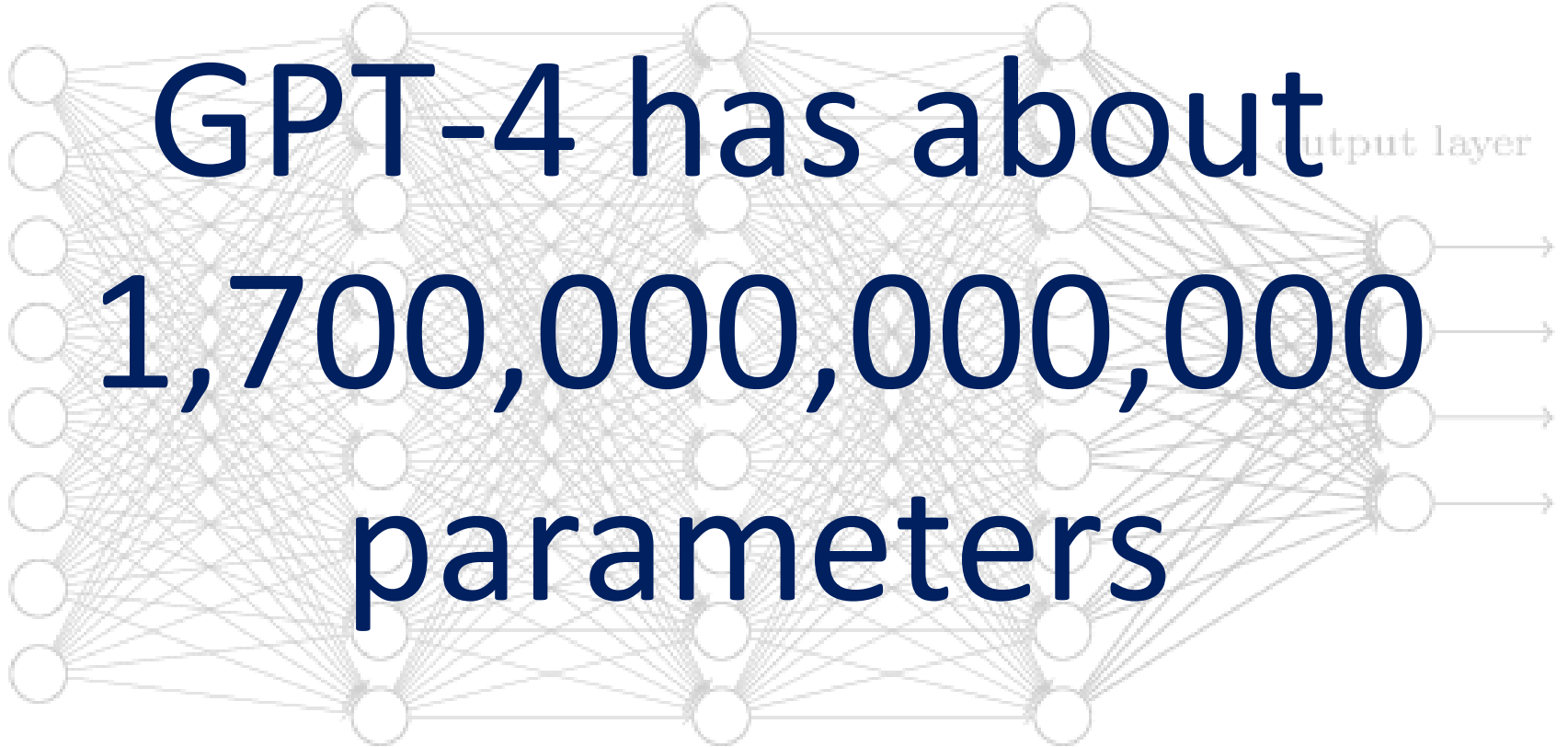input layer    hidden layer 1    hidden layer 2    hidden layer 3    output layer

## Small models (<= 100b parameters)

| ELMo | GPT-1 | BERT | RoBERTa | Transformer ELMo | GPT-2 | Megatron-LM | LLaMA | Chinchilla | YaLM | ERNIE |
|------|-------|------|---------|------------------|-------|-------------|-------|------------|------|-------|
| 94M | 117M | 340M | 354M | 465M | 1.5B | 8.3B | 65B | 80B | 100B | 100B |
| Ai2 | OpenAI | Google | Meta | Ai2 | OpenAI | nVIDIA | Meta | DeepMind | Yandex | Baidu 百度 |

## Large models (>100b parameters)

The base of ChatGPT

| LaMDA | GPT-3 | Jurassic-1 | Gopher | MT-NLG | PaLM | PaLM-E | GPT-4 |
|-------|-------|------------|--------|--------|------|--------|-------|
| 137B | 175B | 178B | 280B | 530B | 540B | 562B | |
| Google | OpenAI | AI21 labs | DeepMind | nVIDIA | Google | Google | OpenAI |

Parent Google

Undisclosed number of parameters
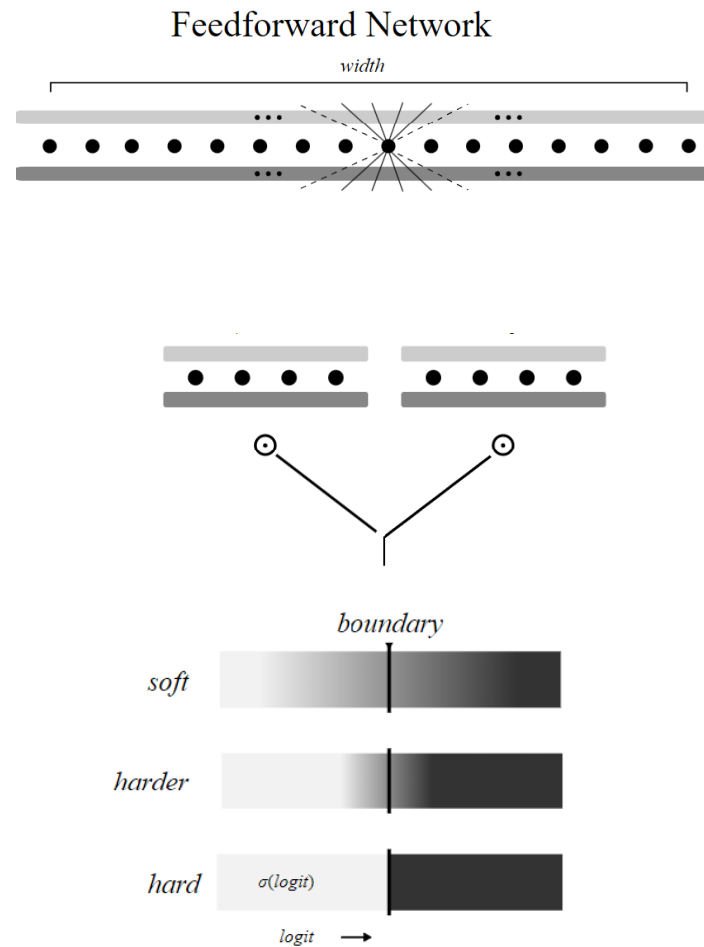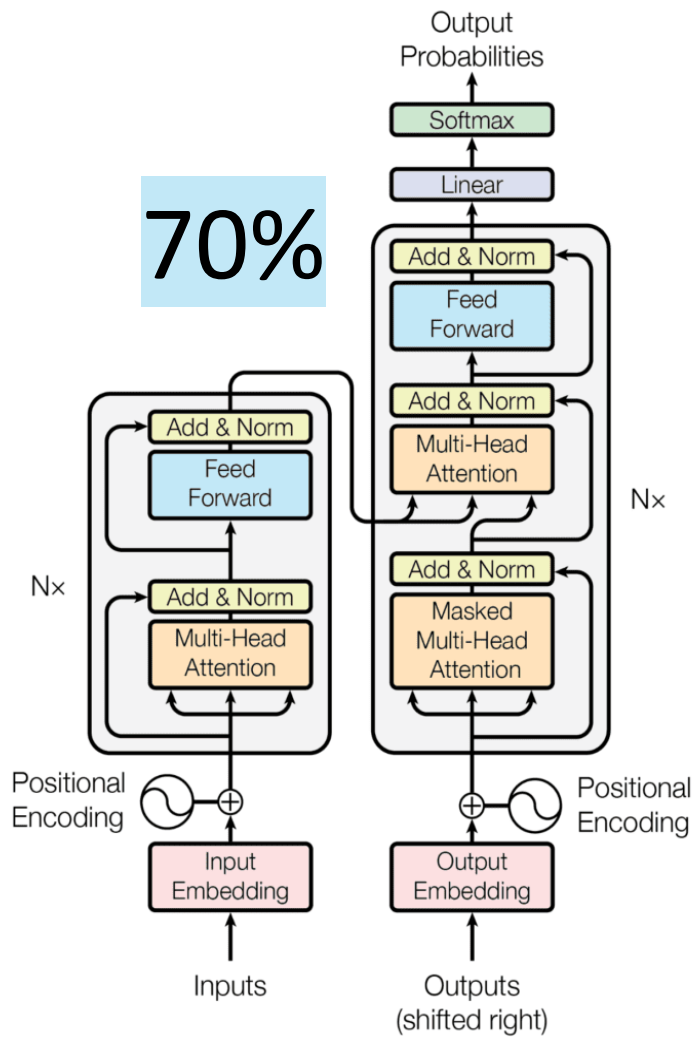
70%

# Exponentially Faster Language Modeling

**Peter Belcak and Roger Wattenhofer**
ETH Zürich
{`belcak,wattenhofer`}@ethz.ch

| Model | $N_T$ | $N_I/N_T$ | RTE | MRPC | STSB | SST-2 | MNLI | QNLI | QQP | Avg | CoLA | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| UltraFastBERT-1x11-long | 4095 | 0.3% | 60.7 | 87.5 | 86.4 | 89.9 | 81.3 | 89.7 | 87.6 | 83.0 | 35.1 | 77.7 |
| **External Baselines** | | | | | | | | | | | | |
| OpenAI GPT | 3072 | 100% | 56.0 | 82.3 | 80.0 | 91.3 | 81.4 | 87.4 | 70.3 | 78.8 | 45.4 | 75.1 |
| DistilBERT | 3072 | 100% | 59.9 | 87.5 | 86.9 | 91.3 | 82.2 | 89.2 | 71.3 | 81.2 | 52.1 | 77.6 |
| BERT-base | 3072 | 100% | 66.4 | 88.9 | 85.8 | 93.5 | 83.4 | 90.5 | 71.2 | 83.0 | 51.3 | 79.6 |

▲ Exponentially faster language modelling (arxiv.org)
300 points by born-jre 10 days ago | hide | past | favorite | 137 comments

add comment

**NEW** Get trending papers in your email inbox once a day!  Subscribe

# Daily Papers
by 🌍 AK

Search by arxiv id or title

**NOV**
**21**
◄ ►

▲ WithinReason 10 days ago | next [–]
Link to previous paper:

▲ qntty 10 days ago | prev | next [–]

## According to scientists, we only use 0.3% of our neural networks. Imagine if we could use 100%.

r/M...

Beiträge

Gepostet von  u/lexected vor 11 Tagen

This study introduces the "Ultra Fast BERT" model, designed to theoretically enhance inference speed by up to 341 times. The increase in speed is possible by replacing the standard FeedForward layers within the Attention mechanism.

1/4

## [R] Exponentially Faster Language Modelling

Research

181

**TL;DR:** Organize your neurons into a tree to get 78x faster inference (theoretical limit is 341x).

This was demonstrated on BERT-base, where this change preserved 96% of its downstream GLUE performance. For a quick comparison, DistilBERT offers 1.6x acceleration while preserving 97% of GLUE performance.

This is a HuggingFace Featured Paper from 11/21/2023.

Paper: https://arxiv.org/abs/2311.10770

Code: https://github.com/pbelcak/UltraFastBERT

| Exponentially Faster Language Modeling | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | $N_T$ | $N_I/N_T$ | RTE | MRPC | STSB | SST-2 | MNLI | QNLI | QQP | Avg | CoLA | Avg |
| **Baselines** | | | | | | | | | | | | |
| crammedBERT-3072 | 4095 | 100.0% | 58.8 | 87.6 | 85.2 | 91.9 | 82.8 | 90.4 | 89.0 | 83.6 | 45.0 | 79.3 |
| crammedBERT-4095 | 3072 | 100.0% | 57.6 | 89.1 | 85.9 | 91.9 | 81.3 | 90.9 | 87.6 | 83.2 | 47.9 | 79.3 |
| **UltraFastBERTs** | | | | | | | | | | | | |

BRAIN SYSTEMS

**1**
- Intuition
- Involuntary control
- Effortlessness
- Innate skills
- Speed
- Reflex

**2**
- Patience
- Concentration
- Constructed thoughts
- Analytics
- Effort
- Reason

**Classification**

**Regression**

**Forecast**

**Clustering**

**Machine Translation**

**Computer Vision**
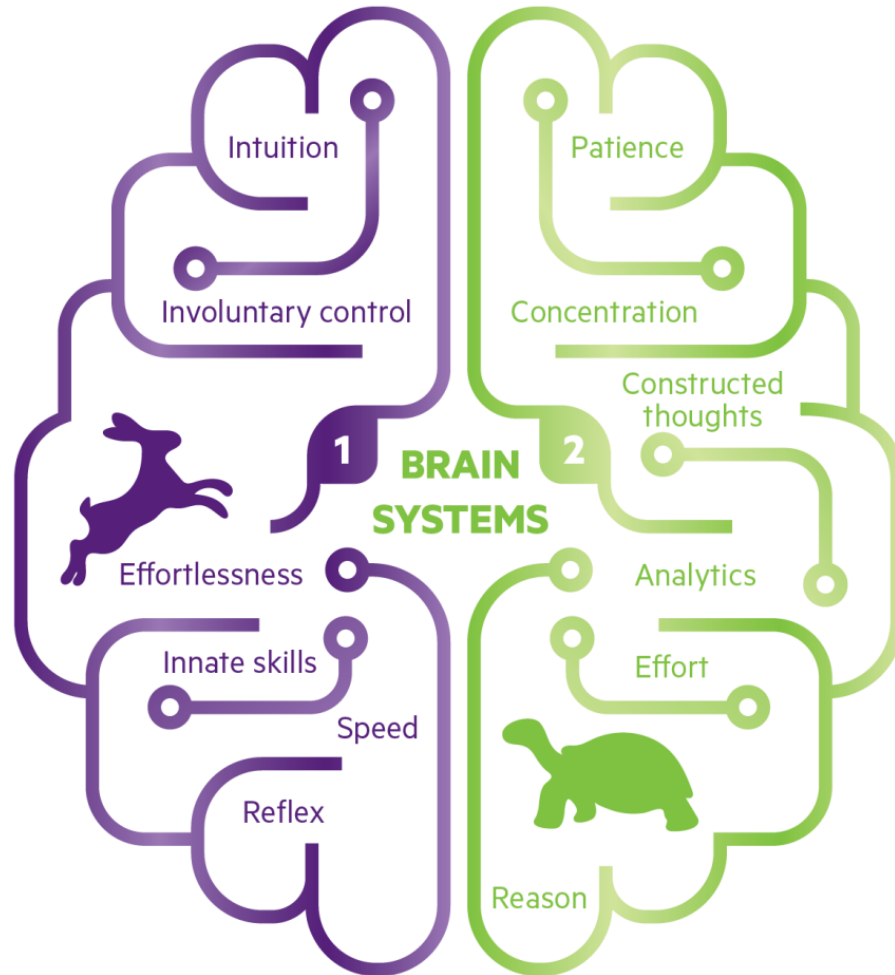
**Generative Art**

**Reinforcement learning**

**uter**
**on**

**Generative Art**

**Reinfor**
**lear**

Refik Anadol, Unsupervised, MOMA

**BRAIN SYSTEMS**

**1**
- Intuition
- Involuntary control
- Effortlessness
- Innate skills
- Speed
- Reflex

**2**
- Patience
- Concentration
- Constructed thoughts
- Analytics
- Effort
- Reason

**User**

what is 5 times 7?

**User**

what is 384903720 times 538982734?

**Assistant**

5 times 7 equals 35.

**Assistant**

384,903,720 times 538,982,734 equals 207,559,081,508,668,480.

Calculater
575-35

384903720 × 538982734 =

207,456,459,332,370,480

314

159 x
_____
28²2³6

157²00

31400
_____
49926

**User**

what is 384903720 times 538982734?

**Assistant**

384,903,720 times 538,982,734 equals 207,559,081,508,668,480.

# Thinking Slow Benchmark?

# FACT: Learning Governing Abstractions Behind Integer Sequences

Peter Belcák, Ard Kastrati, Flavio Schenker, Roger Wattenhofer

| Task | Input | Output | Within class | Across classes | | | |
|------|-------|--------|------|------|------|------|------|
| | | | | | Class: Not Given | | |
| | | | | Class: Polynom | Class: Not Given | | |
| Sequence Unmasking | 1, 2, □, 4, 5, □, 7 | 1, 2, 3, 4, 5, 6, 7 | | Class: Polynom | Class: Not Given | | |
| Sequence Contination | 1, 1, 2, 3, 5, 8, ? | 13 | | Class: Polynom | | Class: Not Given | |
| Sequence Similarity | 1, 2, 3, 4, 5 ... 2, 4, 6, 8, 10 … | Similar | | Is it periodic? | | | |
| Sequence Classif. | 0, 1, 2, 0, 1, 2… | Periodic | | | | | |

数独通信
数独好きがつながる雑誌

季刊 パズル通信 ニコリ
2022年 春号
特集 遊園地は良いぞ
Vol.178

別冊
…イアント
？

ひらめきパズル
ニコリの
謎解き

with
クロスワード
4

漢字
パズル
コレクション

絵むすび
コレクショ…

…解いて気持ち良く、…健康と平和ょ届け！

も数独伝道師!?
…の「教え方」
…数独の正解者

…いパズルや、読者参加…
…シャツなどが当たる懸…

解けば
…ひらめ…

# Simon Tatham's Portable Puzzle Collection



**Simon Tatham's Puzzles** 4+
Greg Hewgill
Designed for iPad

★★★★★ 4.8 • 171 Ratings

Free



# Simon Tatham's Puzzles

**Chris Boyle**

| 4.8★ | 500K+ | E |
|---|---|---|
| 14.5K reviews | Downloads | Everyone ⓘ |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Black Box | Bridges | Cube | Dominosa | Fifteen | Filling | Flip | Flood |
| Galaxies | Guess | Inertia | Keen | Light Up | Loopy | Magnets | Map |
| Mines | Mosaic | Net | Netslide | Palisade | Pattern | Pearl | Pegs |
| Range | Rectangles | Same Game | Signpost | Singles | Sixteen | Slant | Solo |
| Tents | Towers | Tracks | Twiddle | Undead | Unequal | Unruly | Untangle |

# Loopy

# Loopy (Takegaki, Slitherlink, Ouroboros, Suriza, …)

# RLP: A Reinforcement Learning Benchmark for Neural Algorithmic Reasoning

| Puzzle | Parameters | PPO | DreamerV3 |
|---|---|---|---|
| Netslide | 2x3b1 | $35.3 \pm 0.7$ (100.0%) | $12.0 \pm 0.4$ (100.0%) |
| | 3x3b1 | $4742.1 \pm 2960.1$ (9.2%) | $3586.5 \pm 676.9$ (22.4%) |
| Same Game | 2x3c3s2 | $11.5 \pm 0.1$ (100.0%) | $7.3 \pm 0.2$ (100.0%) |
| | 5x5c3s2 | $1009.3 \pm 1089.4$ (30.5%) | $527.0 \pm 162.0$ (30.2%) |
| Untangle | 4 | $34.9 \pm 10.8$ (100.0%) | $6.3 \pm 0.4$ (100.0%) |
| | 6 | $2294.7 \pm 2121.2$ (96.2%) | $1683.3 \pm 73.7$ (82.0%) |

# Sudoku

| | 4 | 6 | | | 2 | | | |
|---|---|---|---|---|---|---|---|---|
| | 3 | | | | | 7 | | |
| 7 | | 2 | | 9 | 8 | | | |
| | | 5 | | | | | 2 | |
| | 8 | | 5 | | 6 | | 3 | |
| | 2 | | | | | 5 | | |
| | | | 7 | 1 | | 2 | 4 | 6 |
| | | 9 | | | | | 5 | |
| | | | 4 | | | 8 | 9 | |

Sudoku RecGNN (Iterative Solving)
Step 0

Kakuro GNN
Step: 0

# Tents

# Atari Games

# Puzzle Coll.

## BRAIN SYSTEMS

**1**
- Intuition
- Involuntary control
- Effortlessness
- Innate skills
- Speed
- Reflex

**2**
- Patience
- Concentration
- Constructed thoughts
- Analytics
- Effort
- Reason

THOMAS H. CORMEN

CHARLES E. LEISERSON

RONALD L. RIVEST

CLIFFORD STEIN

Over
**1 MILLION**
copies sold
worldwide

INTRODUCTION TO
# ALGORITHMS

**FOURTH EDITION**

# The CLRS Algorithmic Reasoning Benchmark

**Petar Veličković** [1]  **Adrià Puigdomènech Badia** [1]  **David Budden** [1]
**Razvan Pascanu** [1]  **Andrea Banino** [1]  **Misha Dashevskiy** [1]  **Raia Hadsell** [1]  **Charles Blundell** [1]

# SALSA-CLRS: A Sparse and Scalable Benchmark for Algorithmic Reasoning

# ARC
## Abstraction & Reasoning Corpus

On the Measure of Intelligence

François Chollet *
Google,

# Abstract Visual Reasoning Enabled by Language

Giacomo Camposampiero*    Loïc Houmard*    Benjamin Estermann    Joël Mathys

Roger Wattenhofer

ETH Zürich, Switzerland

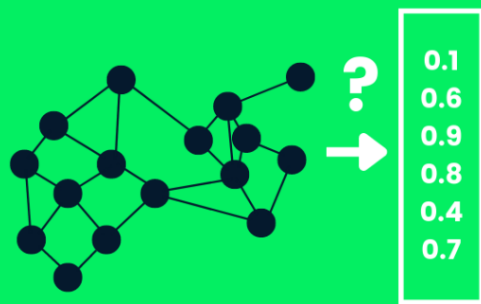{gcamposampie, lhoumard, estermann, jmathys, wattenhofer}@ethz.ch
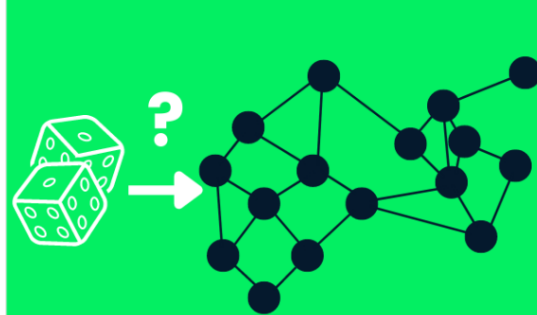
**Graph Classification**

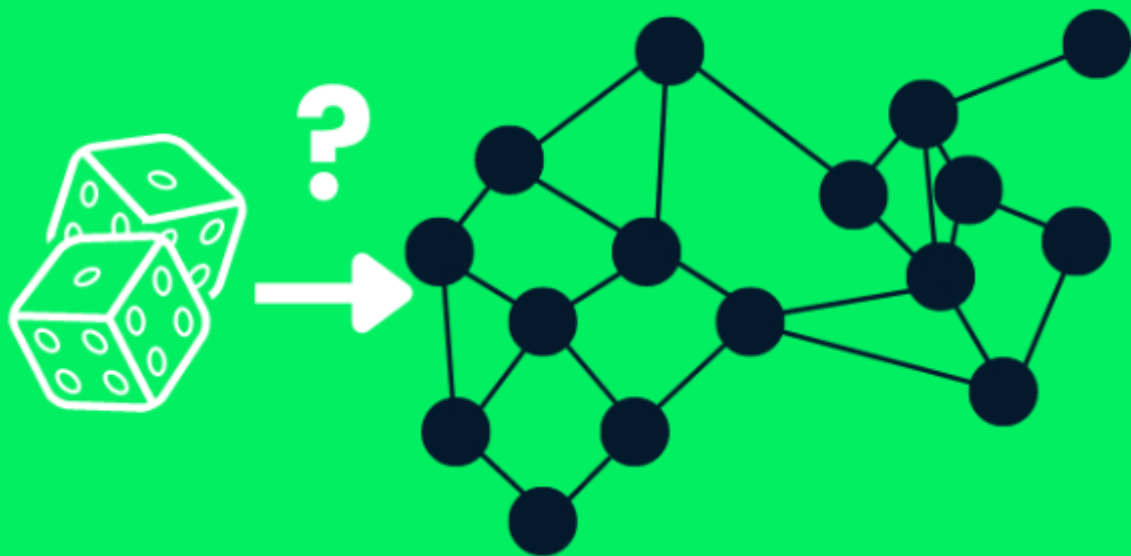**Node Classification**

**Link Prediction**

**Community Detection**

**Graph Embedding**

0.1
0.6
0.9
0.8
0.4
0.7

**Graph Generation**

Graph Generation

# SPECTRE : Spectral Conditioning Helps to Overcome the Expressivity Limits of One-shot Graph Generators

**Karolis Martinkus** [1]   **Andreas Loukas** [* 2]   **Nathanaël Perraudin** [* 3]   **Roger Wattenhofer** [1]

# Discovering Graph Generation Algorithms

**Mihai Babiac, Karolis Martinkus & Roger Wattenhofer**
ETH Zurich
{mbabiac,martinkus,wattenhofer}@ethz.ch

```
1  def outer_loop():
2      for i in range(N):
3          inner_loop()
4
5  def inner_loop():
6      for j in range(i):
7          float00 = random(0, 1)
8          bool00 = float00 < 0.4
9          if bool00:
10             add_edge(i, j)
11
12 outer_loop()
```

```
1  def outer_loop():
2      for i in range(N):
3          int00 = i + n
4          add_edge(i, int00)
5
6          int01 = i % n
7          bool00 = int01 == 0
8          if not bool00:
9              int01 = i + 1
10             add_edge(i, int01)
11
12 outer_loop()
```
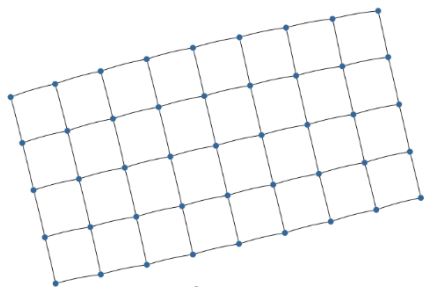
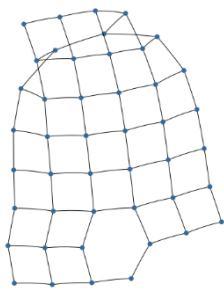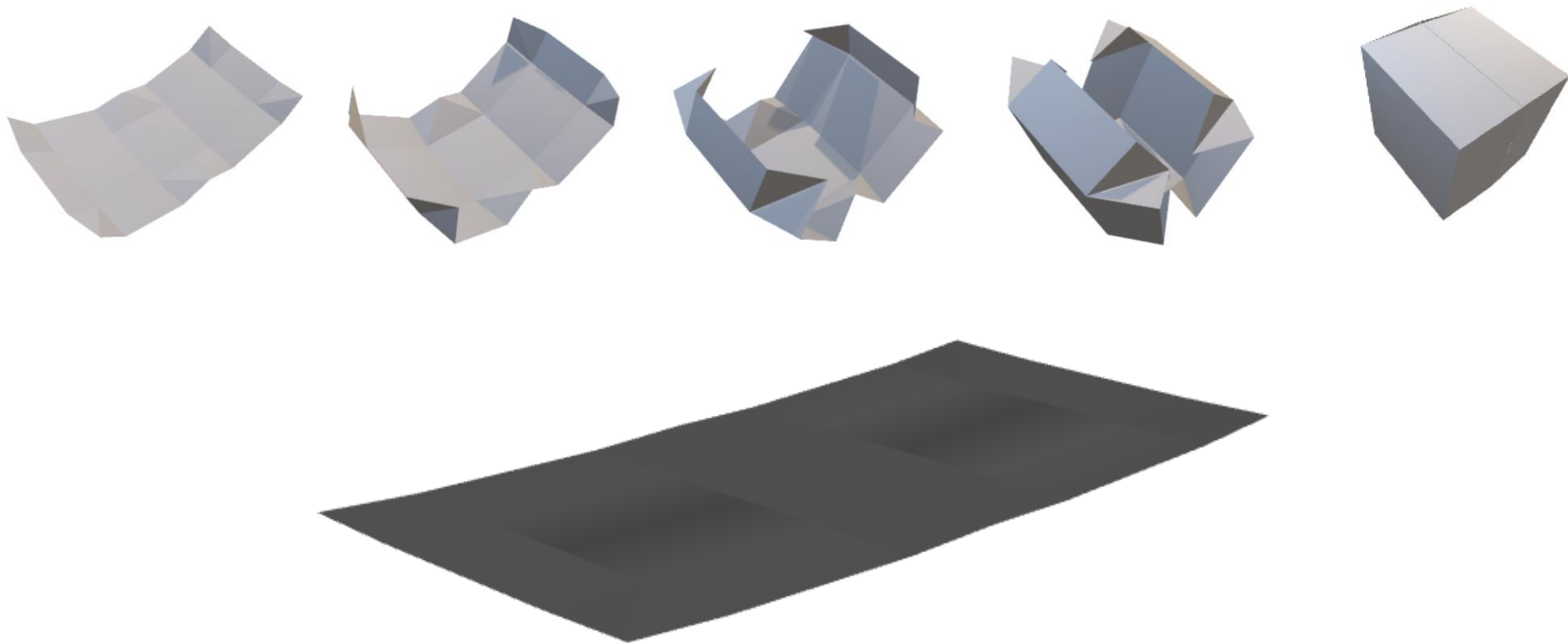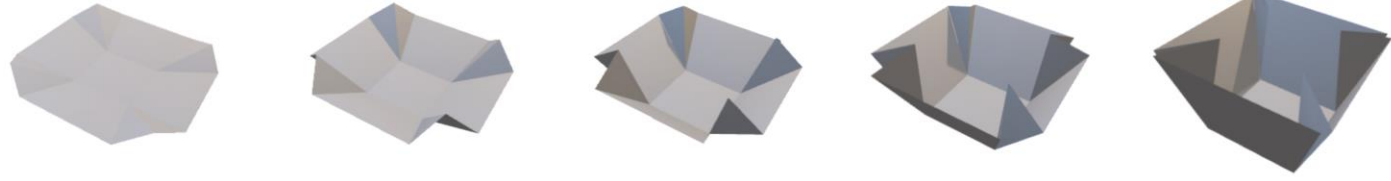|  | Grid | Grid with width | Lobster | Protein |
|---|---|---|---|---|
| Reference | | | | |
| Generated | | | | |

# Automating Rigid Origami Design

Jeremia Geiger, Karolis Martinkus, Oliver Richter, Roger Wattenhofer
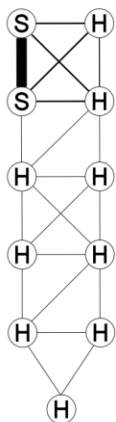
# GraphChef: Learning the Recipe of Your Dataset



(a)  (b)  (c)  (d)  (e)

120"

100"

85"

The Bigger Picture

# Atari Games

# Puzzle Coll.

## BRAIN SYSTEMS

**1**

Intuition
Involuntary control
Effortlessness
Innate skills
Speed
Reflex

**2**

Patience
Concentration
Constructed thoughts
Analytics
Effort
Reason

Tents, Range, Mines, Galaxies, Black Box, Towers, Rectangles, Mosaic, Guess, Bridges, Tracks, Same Game, Net, Inertia, Cube, Twiddle, Signpost, Netslide, Keen, Dominosa, Undead, Singles, Palisade, Light Up, Fifteen, Unequal, Sixteen, Pattern, Loopy, Filling, Unruly, Slant, Pearl, Magnets, Flip, Untangle, Solo, Pegs, Map, Flood
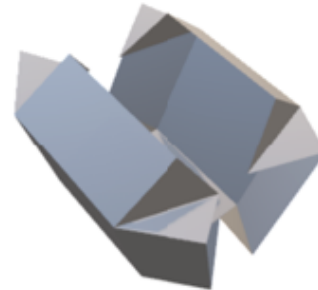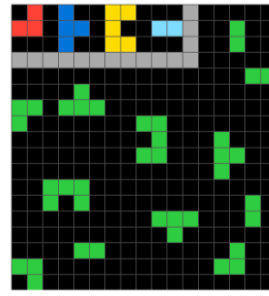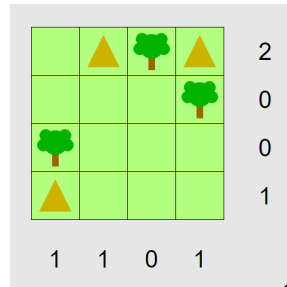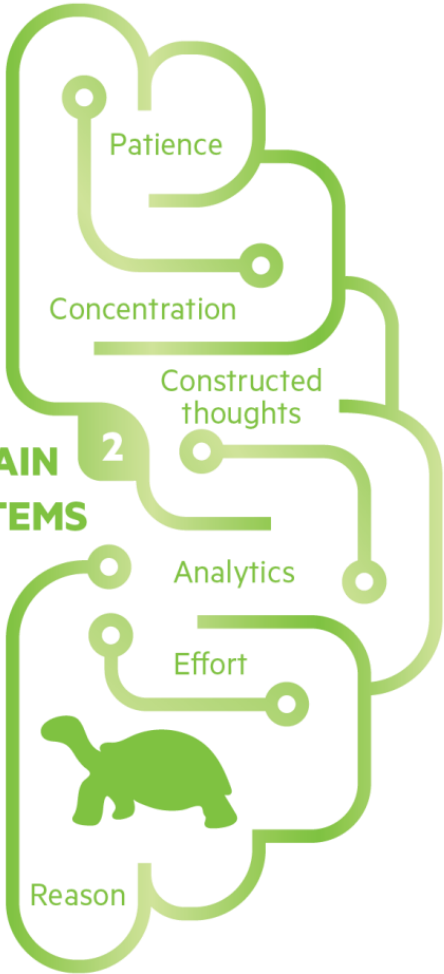
# Problem = Graph

# Solution = Dist. Learning

# Thank You!

## *Any questions or comments?*

*Thanks to co-authors: Peter Belcák, Benjamin Estermann, Lukas Faber, Florian Grötschla, Ard Kastrati, Luca Lanzendörfer, Karolis Martinkus, Joël Mathys, etc.*

*Roger Wattenhofer*