

# Photonic computing for large-scale AI applications

A (non-linear) entrepreneurial journey from optical computing to AI agents

**Laurent DAUDET**

LightOn, Co-Founder

Université Paris Cité, Professor

IT4I Workshop on Quantum Computing, 3 December 2025

# A spinoff from university research

4 Co-founders

A team of ~45



**Igor Carron**

**CEO**

Management of complex nuclear engineering/ aerospace projects

- Former Assistant Director at Spacecraft Technology Center at Texas A&M Univ.



**Laurent Daudet**

**Expert in Signal Processing**

- Physics professor at Université Paris Cité



**Florent Krzakala**

**ML Advisor**

- Professor of Physics and Electrical Engineering at EPFL



**Sylvain Gigan**

**Optics Advisor**

- Physics professor at Sorbonne University



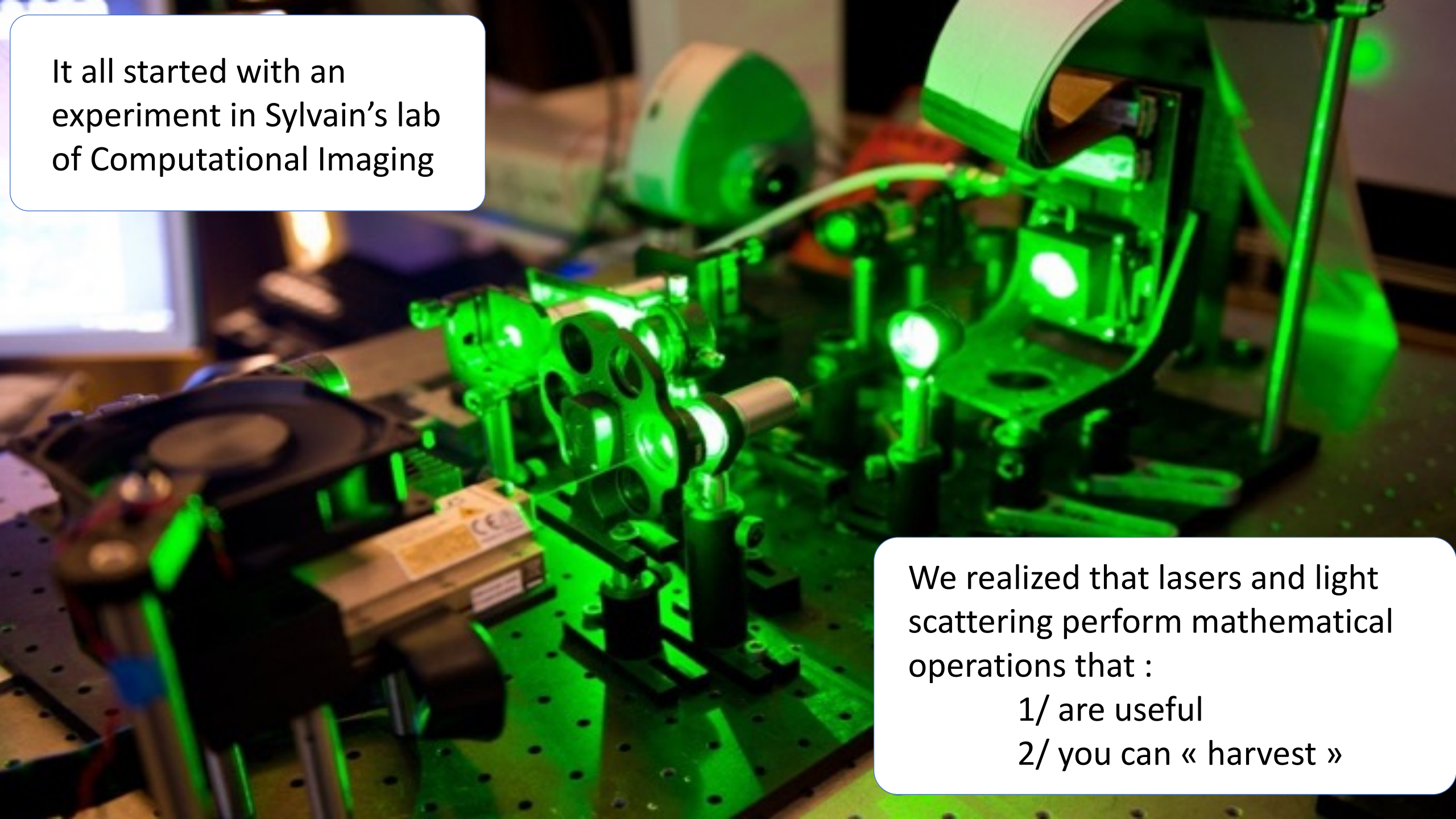
# How LightOn *actually* started

## Procrastination



## Coffee





It all started with an experiment in Sylvain's lab of Computational Imaging

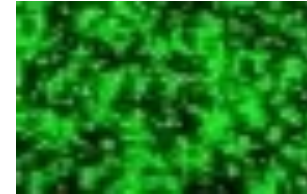
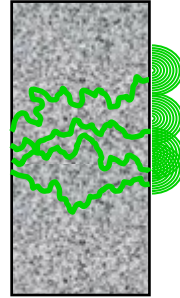
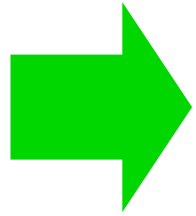
We realized that lasers and light scattering perform mathematical operations that :

1/ are useful

2/ you can « harvest »

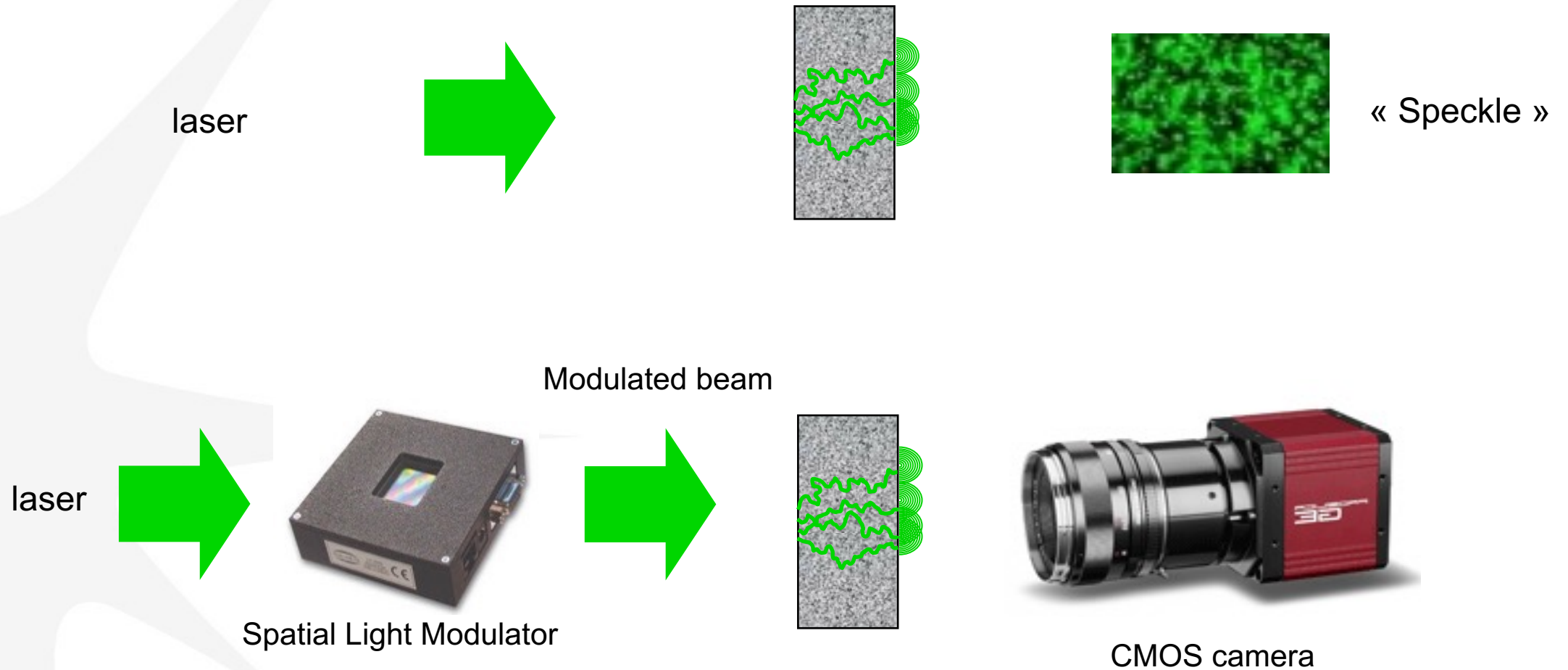
# Scattering: a coherent process

laser

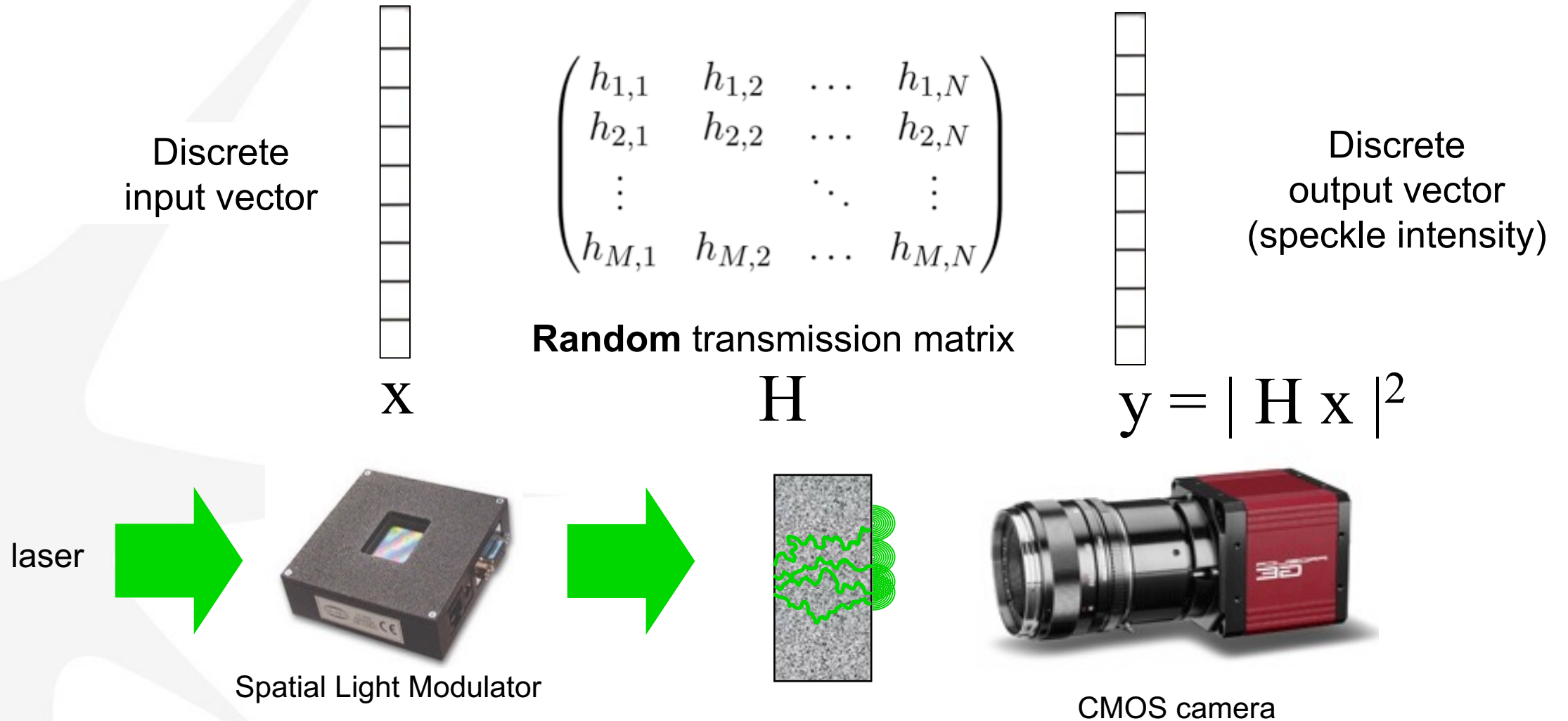


« Speckle »

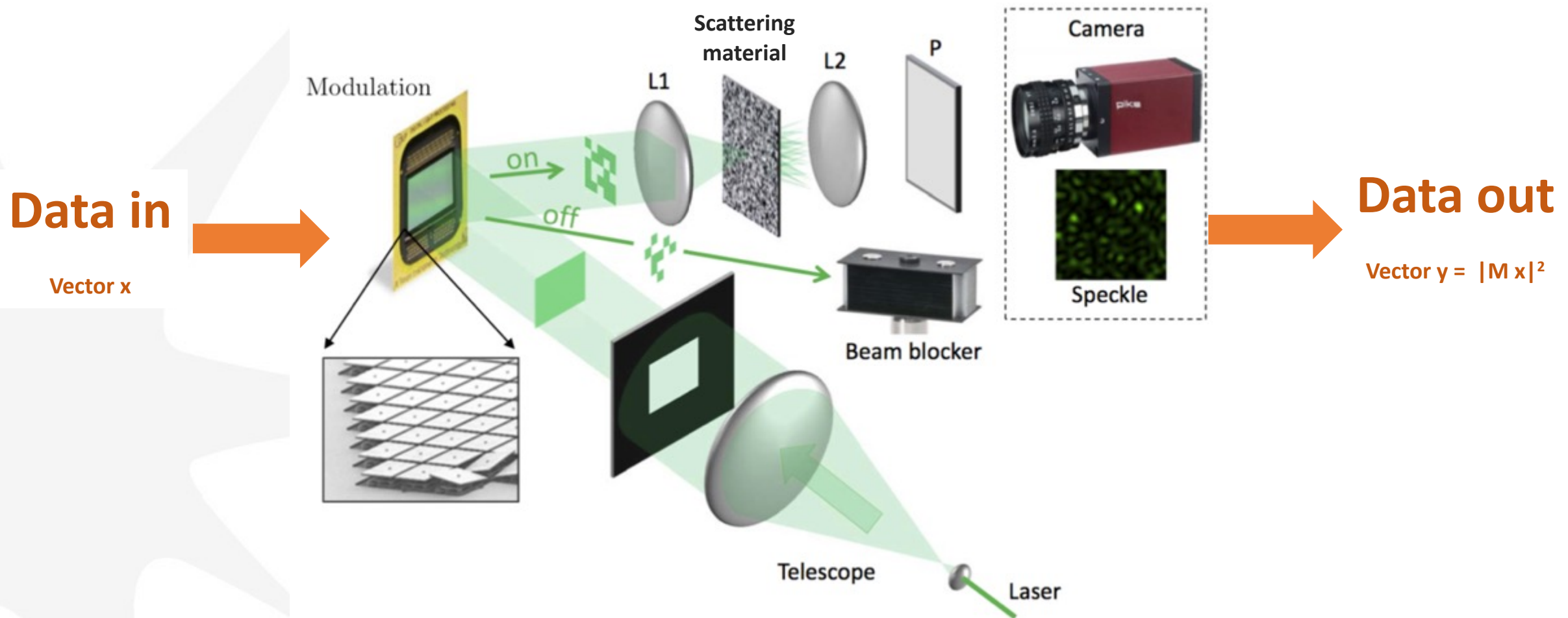
# Scattering: a coherent process



# Scattering: a coherent process

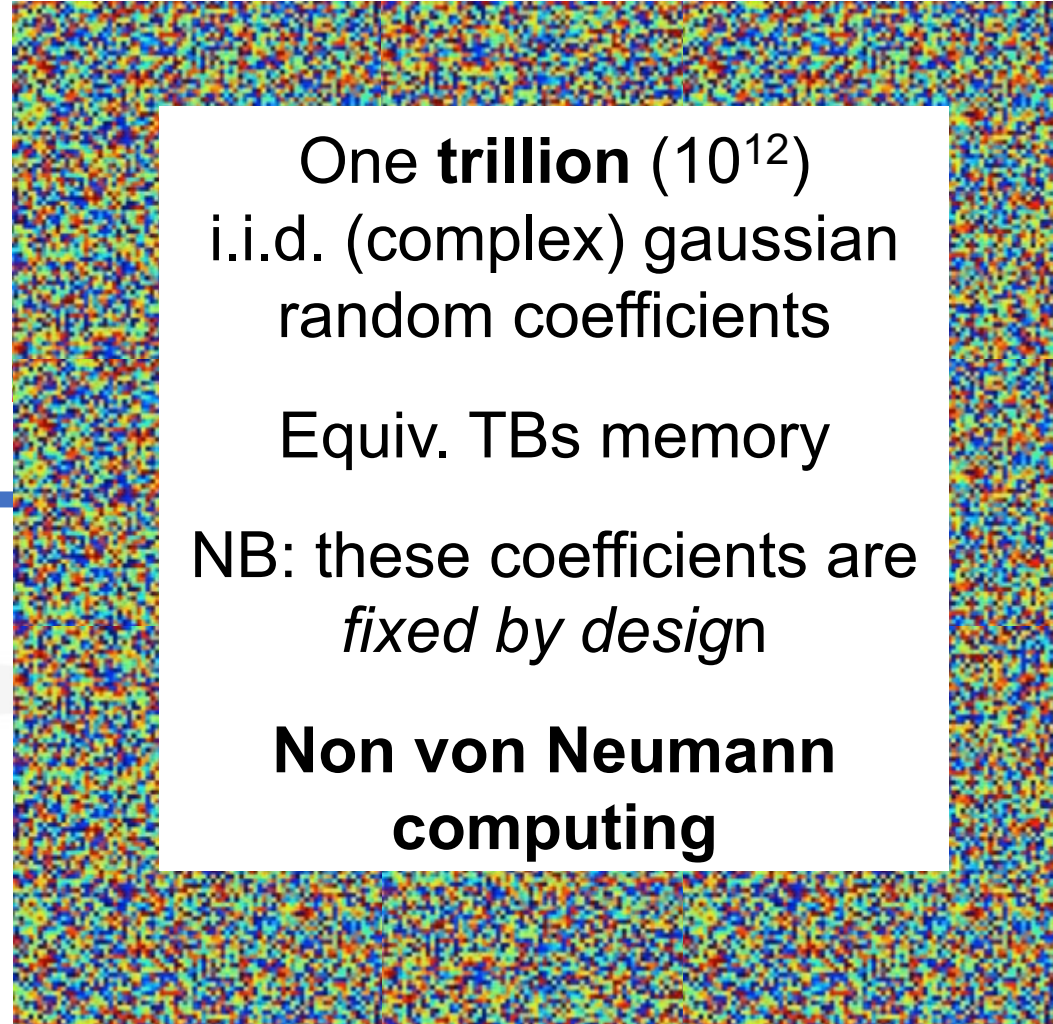


# Optical Processing Unit technology



# Matrix-vector multiplication through light scattering

1 million  
independent  
input  
pixels



One **trillion** ( $10^{12}$ )  
i.i.d. (complex) gaussian  
random coefficients  
Equiv. TBs memory  
NB: these coefficients are  
*fixed by design*  
**Non von Neumann  
computing**



1 million  
independent  
output  
pixels

- Random Projections act as distance-preserving point cloud embeddings

## Johnson-Lindenstrauss Lemma (1984)

**Lemma** For any  $0 < \epsilon < 1$  and any integer  $n$  let  $k$  be a positive integer such that

$$k \geq \frac{24}{3\epsilon^2 - 2\epsilon^3} \log n$$

then for any set  $A$  of  $n$  points  $\in \mathbb{R}^d$  there exists a map  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that for all  $x_i, x_j \in A$

$$(1 - \epsilon) \|x_i - x_j\|^2 \leq \|f(x_i) - f(x_j)\|^2 \leq (1 + \epsilon) \|x_i - x_j\|^2$$



- NeurIPS 2017 Test of Time Award  
“Random Features for Large-scale Kernel Machines”, Rahimi, Recht, 2008

# The TRL scale

## TECHNOLOGY READINESS LEVEL (TRL)

|             |   |   |
|-------------|---|---|
| RESEARCH    | 9 | ACTUAL SYSTEM PROVEN IN OPERATIONAL ENVIRONMENT           |
|             | 8 | SYSTEM COMPLETE AND QUALIFIED                             |
|             | 7 | SYSTEM PROTOTYPE DEMONSTRATION IN OPERATIONAL ENVIRONMENT |
| DEVELOPMENT | 6 | TECHNOLOGY DEMONSTRATED IN RELEVANT ENVIRONMENT           |
|             | 5 | TECHNOLOGY VALIDATED IN RELEVANT ENVIRONMENT              |
|             | 4 | TECHNOLOGY VALIDATED IN LAB                               |
| DEPLOYMENT  | 3 | EXPERIMENTAL PROOF OF CONCEPT                             |
|             | 2 | TECHNOLOGY CONCEPT FORMULATED                             |
|             | 1 | BASIC PRINCIPLES OBSERVED                                 |

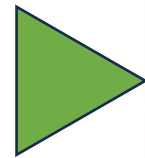
# The first 4 years 2016-2020

## Climbing up the TRL ladder

Lab  
experiment



Component  
Photonic hardware



## TECHNOLOGY READINESS LEVEL (TRL)

|             |   |   |
|-------------|---|---|
| RESEARCH    | 9 | ACTUAL SYSTEM PROVEN IN OPERATIONAL ENVIRONMENT           |
|             | 8 | SYSTEM COMPLETE AND QUALIFIED                             |
|             | 7 | SYSTEM PROTOTYPE DEMONSTRATION IN OPERATIONAL ENVIRONMENT |
| DEVELOPMENT | 6 | TECHNOLOGY DEMONSTRATED IN RELEVANT ENVIRONMENT           |
|             | 5 | TECHNOLOGY VALIDATED IN RELEVANT ENVIRONMENT              |
|             | 4 | TECHNOLOGY VALIDATED IN LAB                               |
| RESEARCH    | 3 | EXPERIMENTAL PROOF OF CONCEPT                             |
|             | 2 | TECHNOLOGY CONCEPT FORMULATED                             |
|             | 1 | BASIC PRINCIPLES OBSERVED                                 |

## LightOn Optical Processing Unit (OPU) the world's first photonic AI co-processor publicly available

**2200 TOPS**

In a single photonic core

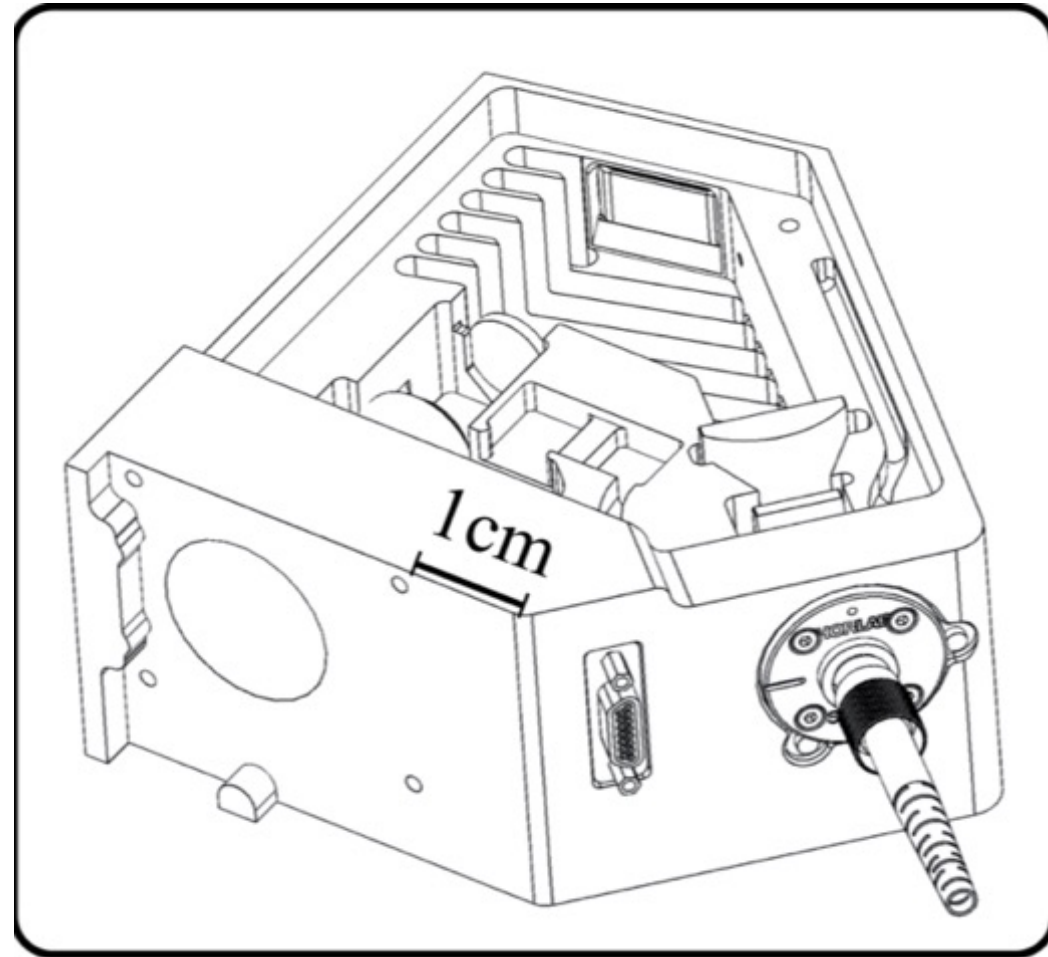
**With only 30 W TDP**

200 times better in #OPS/W than  
NVIDIA V100 GPU boards

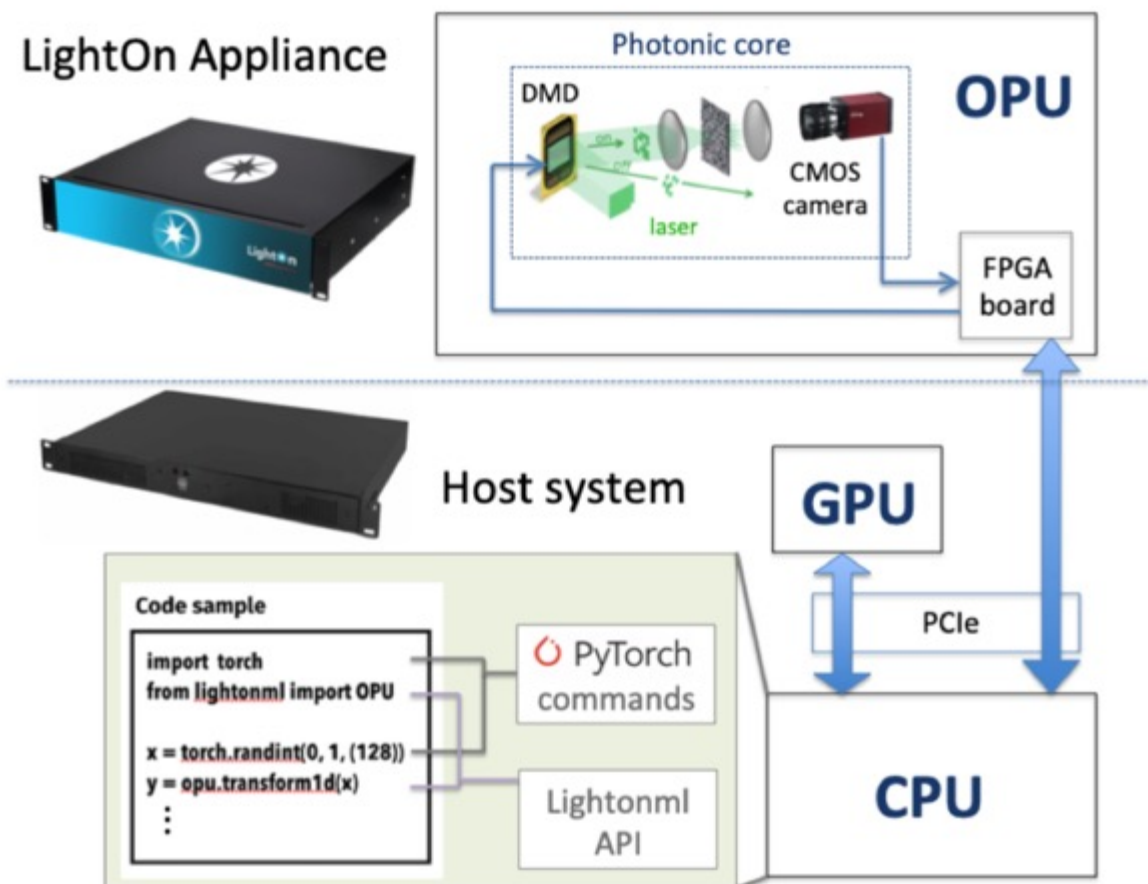


Used by AI researchers through LightOn's cloud platform → build a community

## Optics integration in aluminum monobloc



## Software integration for hybrid data processing architecture



# Developing applications

## Application layer

System

AI software

+

Photonic hardware

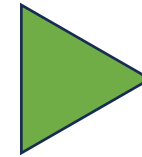


Component

Photonic hardware



Lab  
experiment



# Hybrid computing in AI pipelines



Ivan Dokmanic  
Associate Prof.



NEURAL INFORMATION  
PROCESSING SYSTEMS

NeurIPS 2019



David Rousseau @dhpmrou · 4 avr.

Our talk on analysing #hep data with random matrices in @LightOnIO Optical Processor Unit accepted at @icrep2020 conference ! (Remote in July or Prague beg 2021) #hepml @Laurent\_Daudet @IgorCarron



6



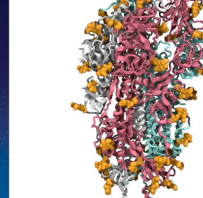
15



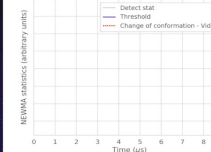
## Analyzing SARS-CoV-2 conformations with a LightOn OPU

### SARS-CoV-2 glycoprotein simulation

initialized in a partially opened state (SVR)



### Changepoint detection (NEWMA-OPU)



LightOn

cloud.lighton.ai

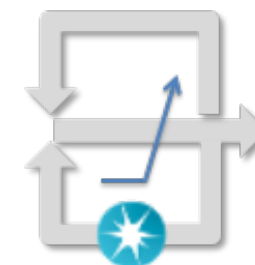
Pre-conditioning  
for RandSVD



Transfer learning  
in Convolutional Deep NN



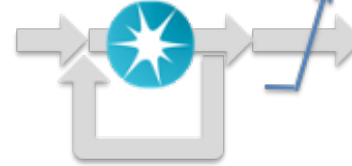
Reinforcement Learning



Sketching  
for changepoint detection



Reservoir computing  
for time / time-space series



Adversarial robustness /  
differential privacy  
by design



Accelerating SARS-CoV2  
Molecular Dynamics Studies  
with Optical Random Features

Amélie  
Chatelain  
LightOn ML  
R&D engineer



NeurIPS 2020

Collaboration  
with Criteo



Collaboration  
with FAIR



# Accelerated scientific computing with Randomized numerical linear algebra

## Randomized Numerical Linear Algebra

*DOE RASC report (Jan 2021):*

*randomized algorithms are "essential to the future of computational science and AI for Science."*

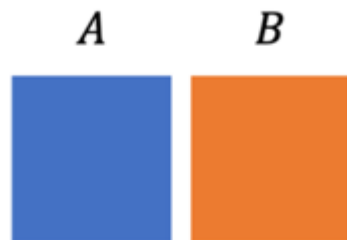
- Approximate matrix multiplications
- Randomized SVD
- ... And much more

### Photonic co-processors in HPC: using LightOn OPU for Randomized Numerical Linear Algebra

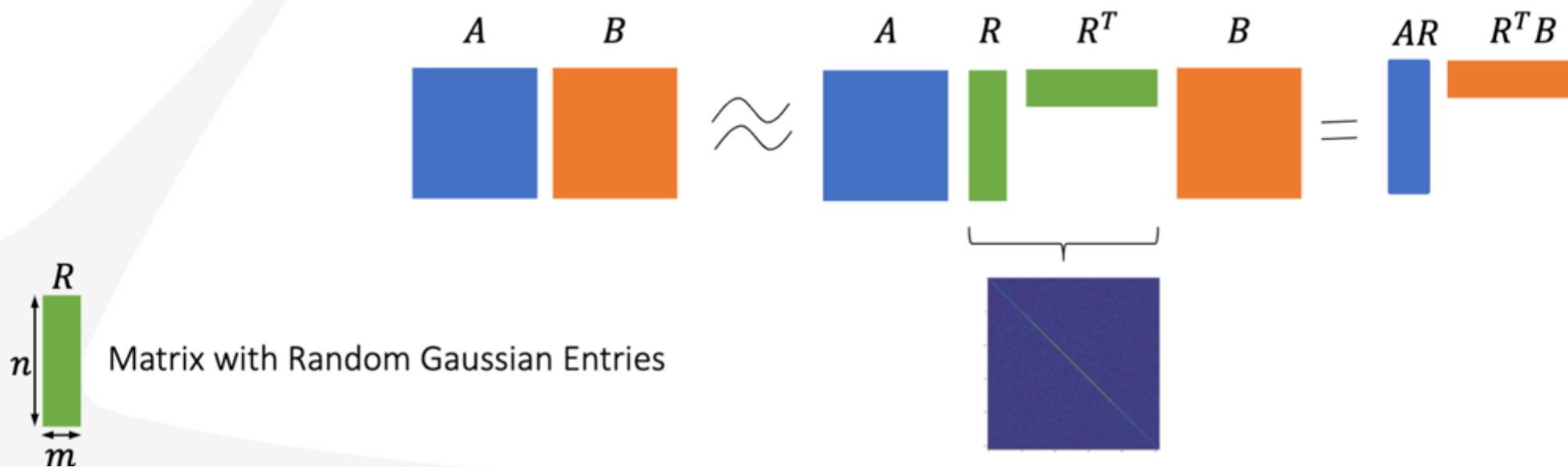
Daniel Hesslow, Alessandro Cappelli, Igor Carron, Laurent Daudet, Raphaël Lafargue,  
Kilian Müller, Ruben Ohana, Gustave Pariente, and Iacopo Poli  
*LightOn. Paris. France.*

See Hesslow D. et al, HotChips2021 proceedings, <https://arxiv.org/abs/2104.14429>

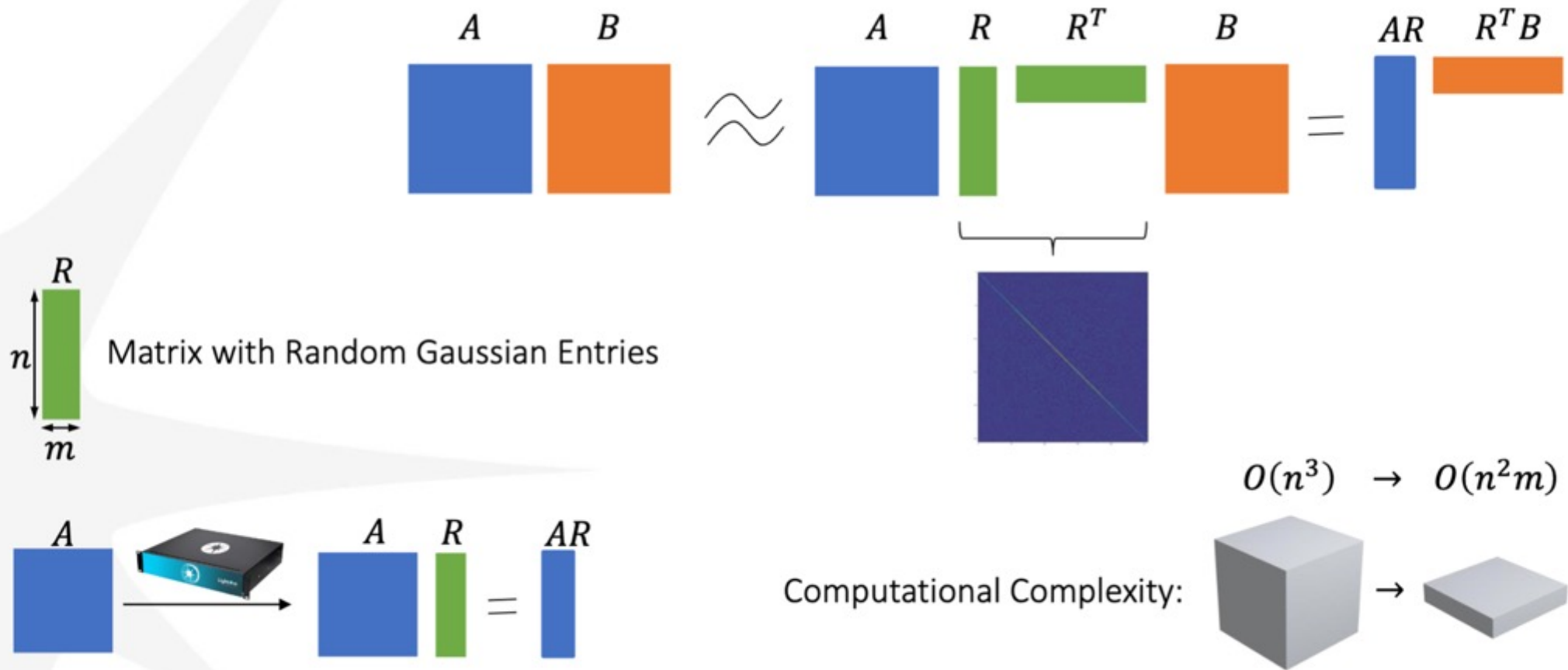
TOY EXAMPLE  
Multiplying 2 matrices



**NB: DON'T DO THIS IN PRACTICE !**



# HPC Use case: Accelerated Scientific Computing

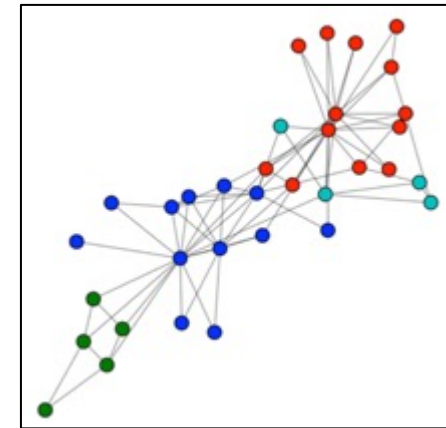


## Randomized Trace estimators – application to Graph Neural Networks

*Hutchinson's trace estimator*

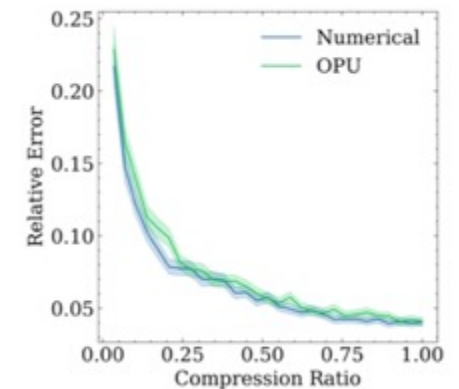
$$\text{Tr}(A) \approx \text{Tr}(RAR^\top)$$

$$\text{Tr}(A^3) \approx \text{Tr}(RA^3R^\top) \approx \text{Tr}((RAR^\top)^3)$$



Community detection in networks  
→ Triangle counting on graphs

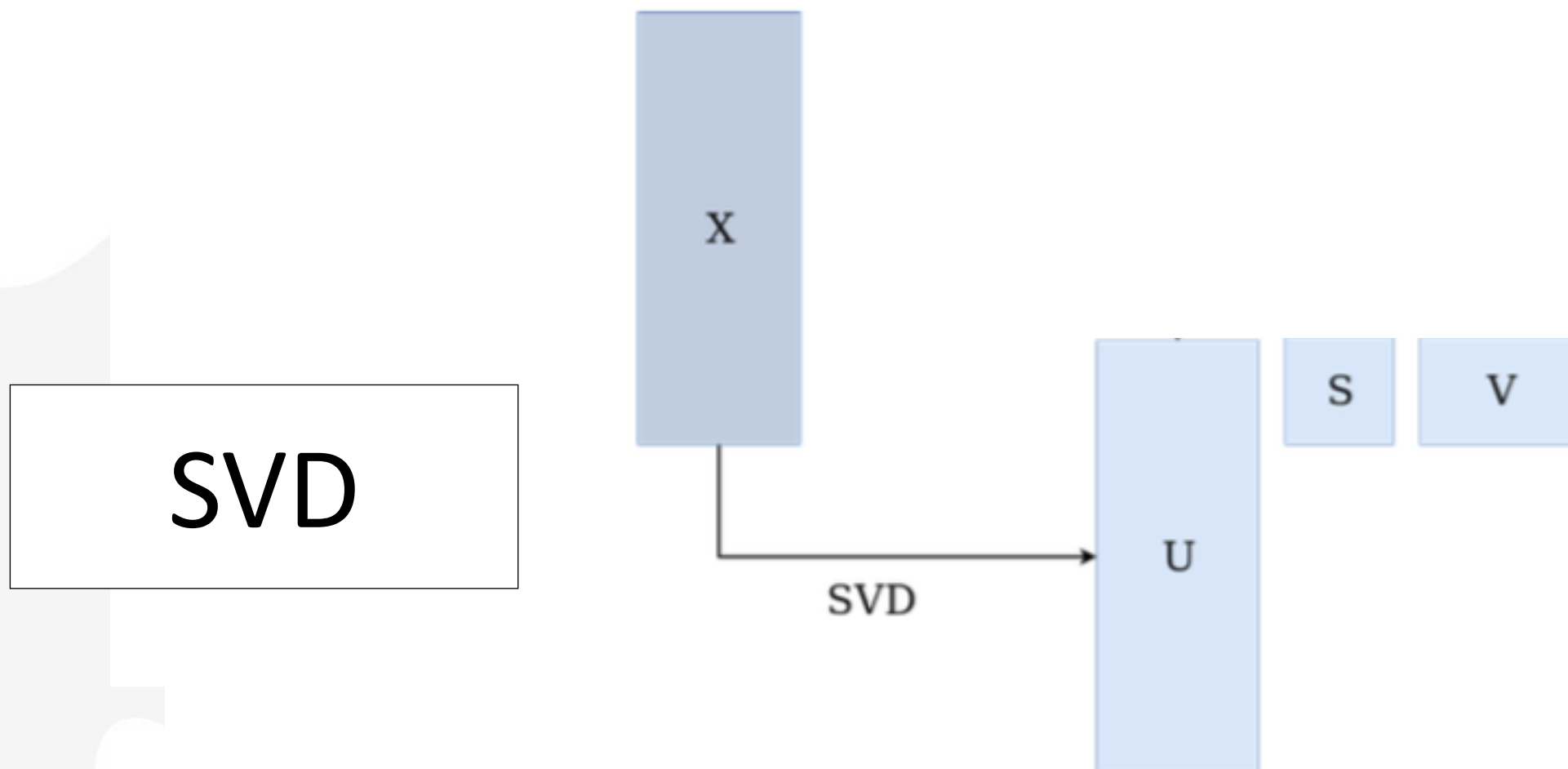
A<sup>3</sup> TRACE ESTIMATION

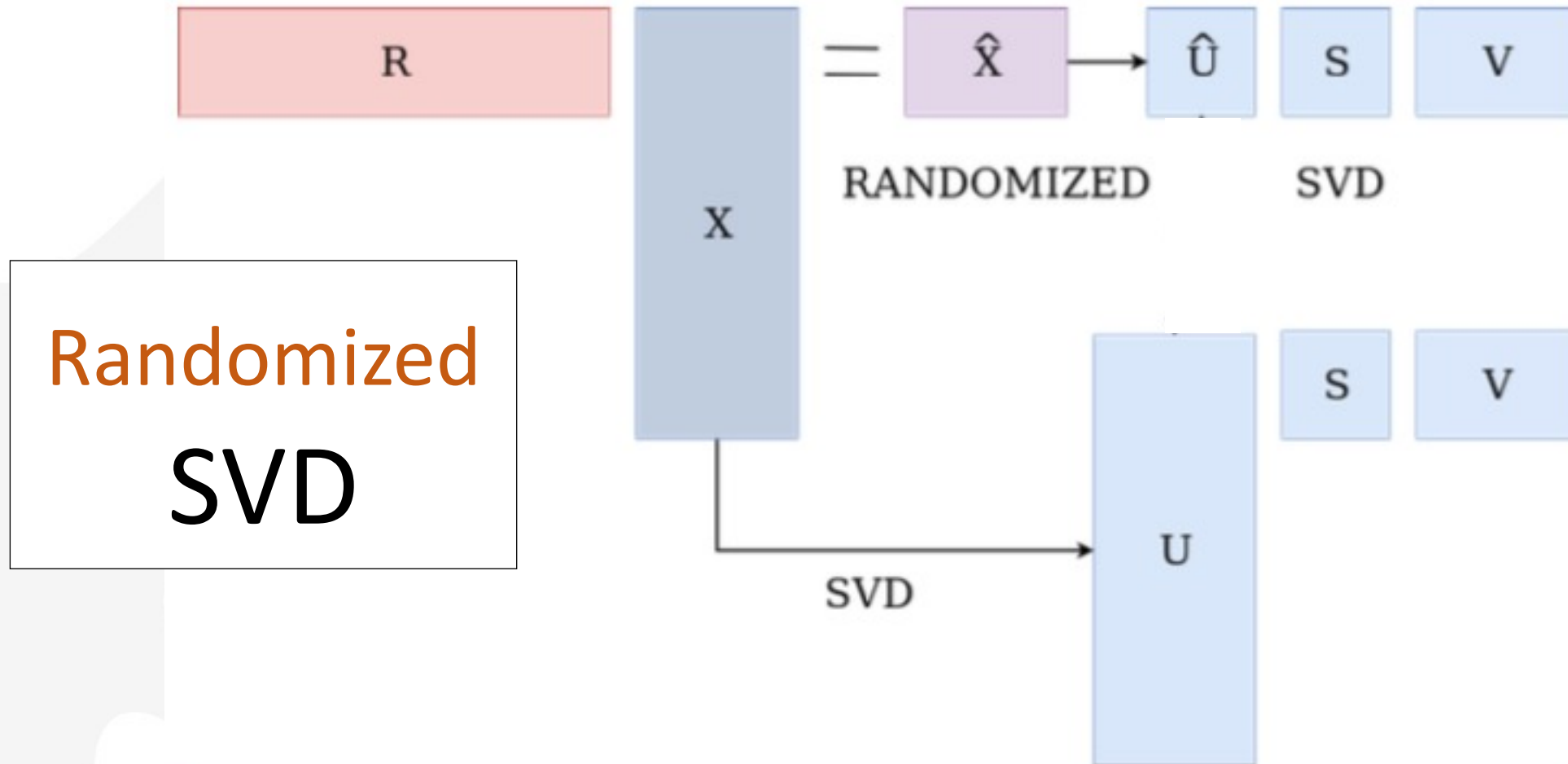


$$O(n^3) \rightarrow O(m^3 + n)$$



(figure from Rossetti et al. Applied Network Science (2019) 4:52)

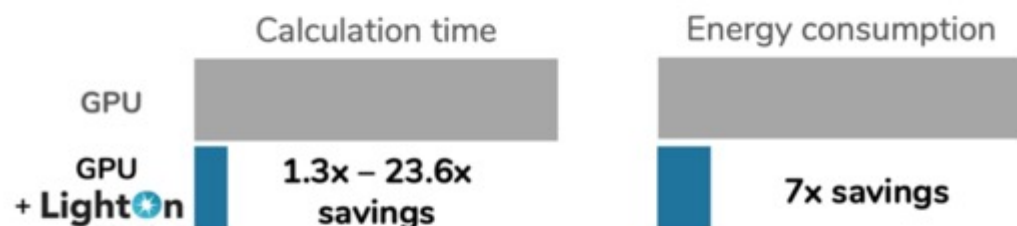




*Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, Halko, N., Martinsson, P., Tropp, J., 2009, arXiv:0909.4061

## LightOn Appliance for AI acceleration

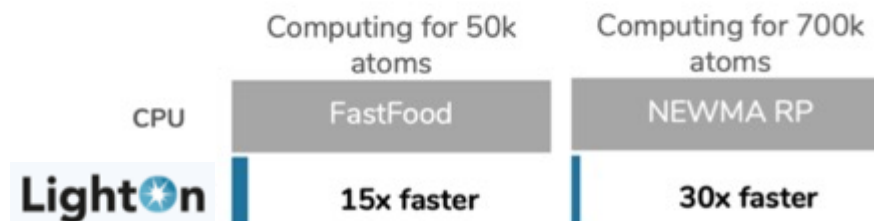
### Approximating Kernel Ridge Regression for classification tasks



Performance gains on Kernel ridge approximation for classification tasks. Dataset qm7 (quantum chemistry), high energy physics, and others. The OPU is compared to an NVIDIA P100 GPU (250 W). GPU RAM limit was hit at 32GB. Results acquired extrapolating to 1M features. OPU: Aurora 1.5 (30 W).

## Real-time AI analysis of large-scale HPC results

### Change detection in Molecular Dynamics

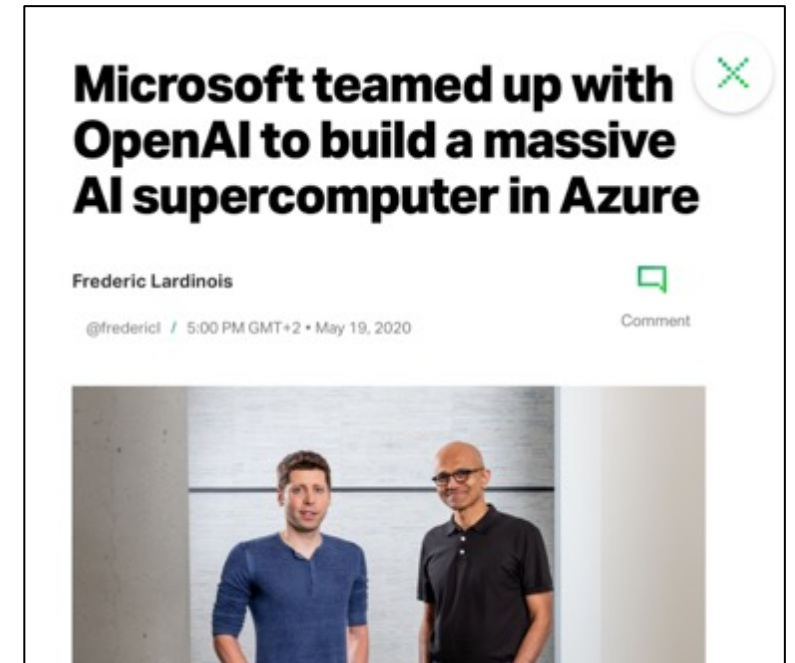


15x faster than FastFood on CPU at 50k atoms. For 700k + atoms, NEWMA RP on OPU is expected to be 30x faster than NEWMA FF on CPU. Library: LightOnML, Dataset: Molecular Dynamics simulations (HPC, Anton), OPU: Aurora 1.5

# The pivot

Training *a single* GPT-3 model :

- 3 Million GPU-hours (on NVIDIA V100s)
- 550 T CO<sub>2</sub> equiv.
- Estimated price 5-10 M \$ for training only



Could OPU technology speed this up ?

# Deploying in operational environments

Lab  
experiment

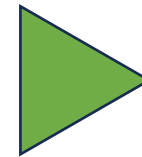


**Component**  
Photonic hardware



**System**

AI software  
+  
Photonic hardware



we put an OPU in a #top500 supercomputer !



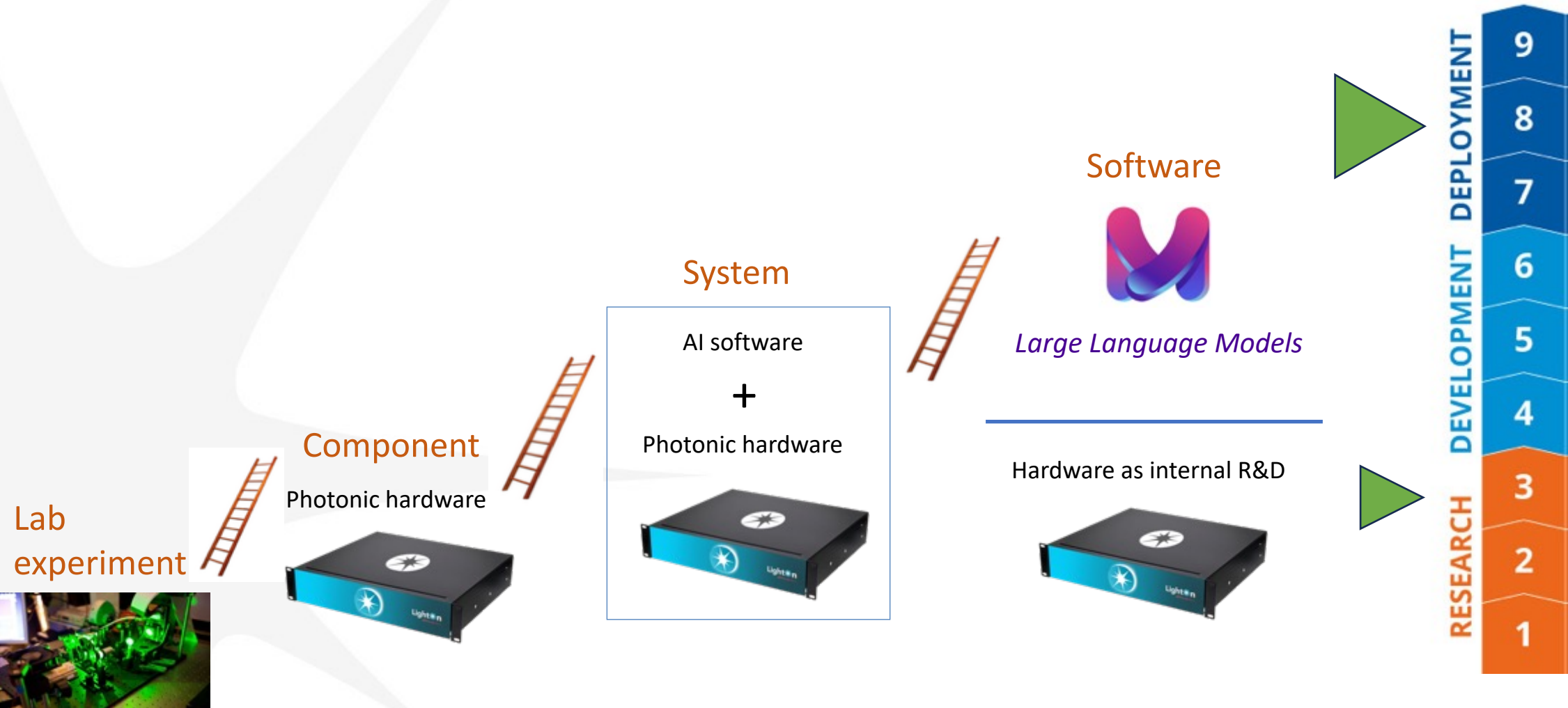




As we first learned how to train LLMs on GPUs ...  
we realized there was an immediate market for it !



# 2021 pivot: decoupling hardware and software



# Contributions to open-source LLMs

Lyra



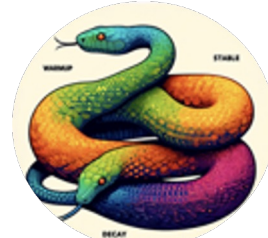
PAGnol



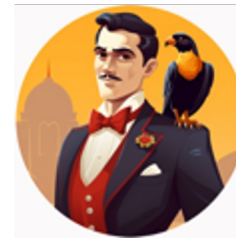
Rita



Mambaoutai



Alfred



ModernBERT



with Answer.AI

# NOOR: world's first LLM in Arabic

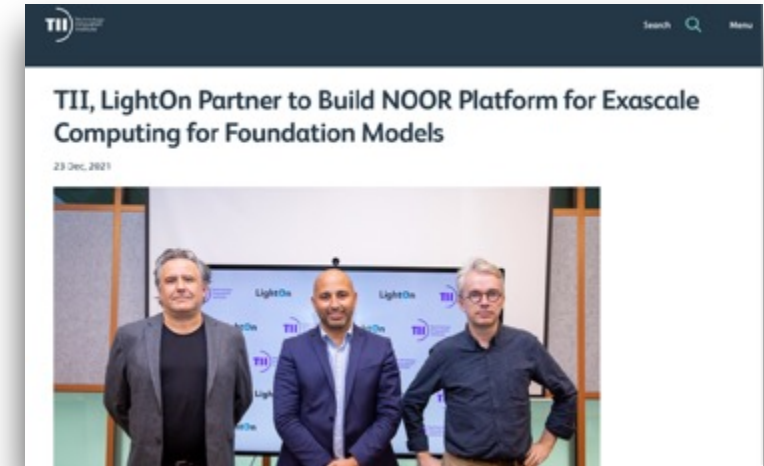


LightOn partnered with TII Technology Innovation Institute Abu Dhabi to build



The world's first Large Language Model in Modern Standard Arabic - released 2022

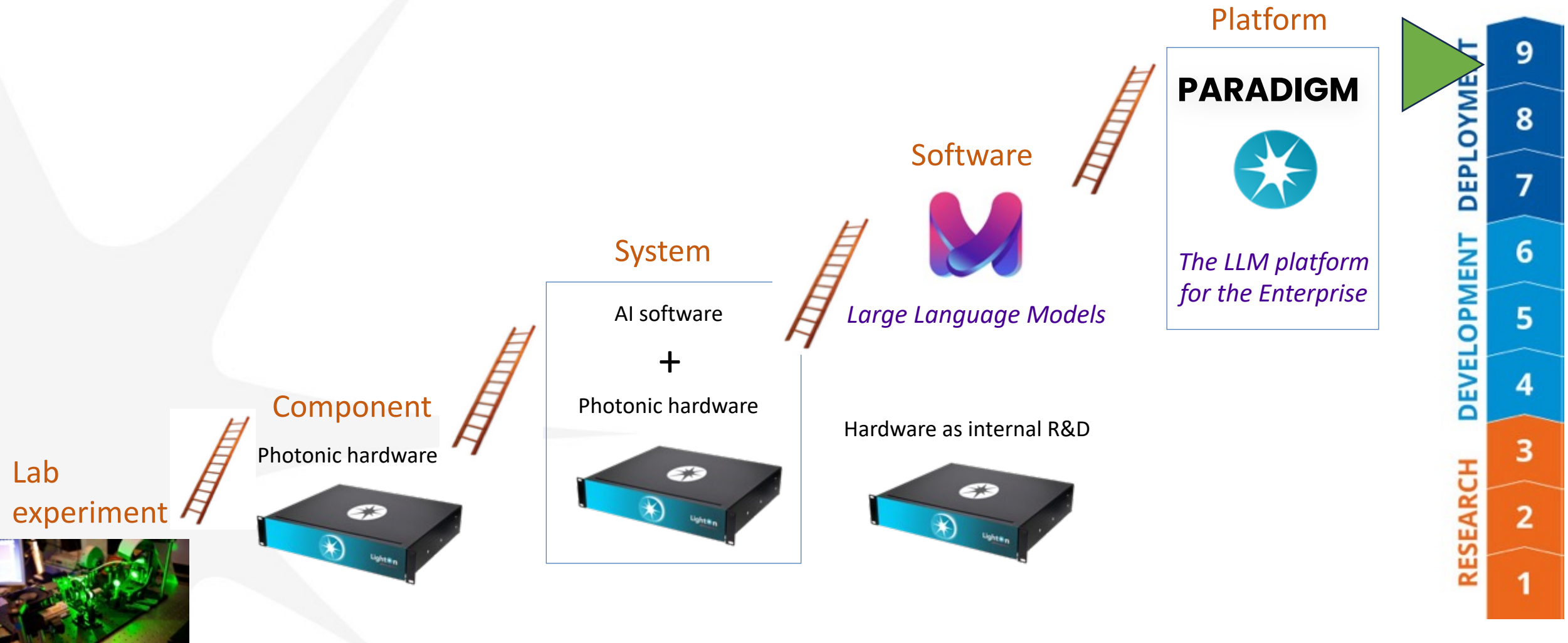
Noor powered the first AI-written journal article in MSA  
AlEtihad News; Nov 19th, 2022



LightOn in 2025

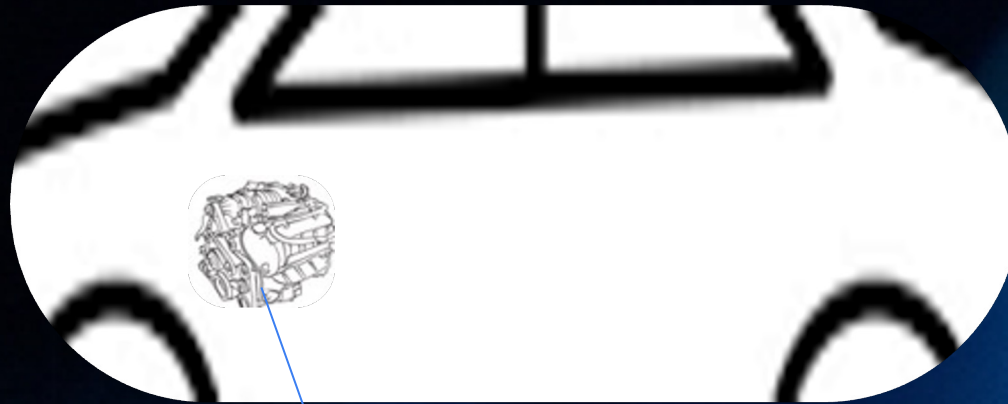
Our mission :  
turn GenAI into an enterprise-ready solution

# Finally reaching the top !



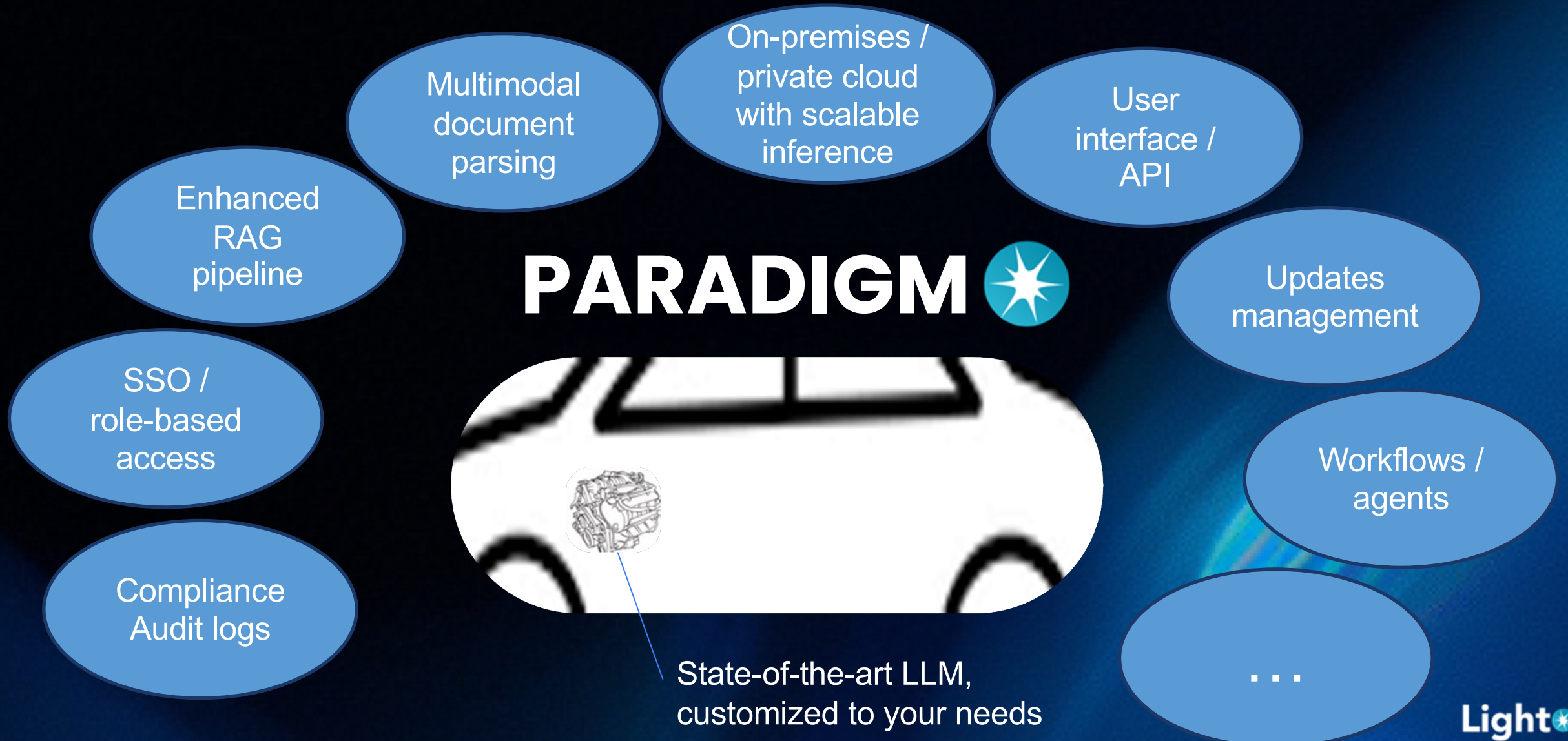
# Turn GenAI into an **enterprise-ready** solution:

- Safe
- Reliable
- Compliant
- Customizable
- Sustainable




State-of-the-art LLM,  
customized to your needs

# Turn GenAI into an enterprise-ready solution



# From documents to insights



**+ NEW CHAT**

Search session

**Today**

AI générative marchés croissance

New Chat

Bonjour, aide concernant LightOn

---

**Analytics**

**Documents**

**Prompts assistant**

**Chat instruction**

---

**Laurent Daudet**

## AI générative marchés croissance

Alfred 4

revenue that is quite significant. The Generative AI market is estimated to generate a revenue of 1.3 trillion dollars, represent about 12% of all technology spending by 2032. This estimation comes from Bloomberg Intelligence, which forecasts the economic potential of generalized Generative AI between 6.1 and 7.9 trillion dollars.



The software segment, where LightOn operates, represents more than half of this Generative AI market by 2032. This means that companies like LightOn, which develop solutions based on language models (LLM) and content generation tools, have a considerable growth potential.

In the French context, it is mentioned that in 2023, only four startups specialized in foundation models were present, including LightOn and MistralAI. This suggests a limited but growing competition in the French market, with international players such as Google, Microsoft, and OpenAI.




Ask your question or type / to access tools

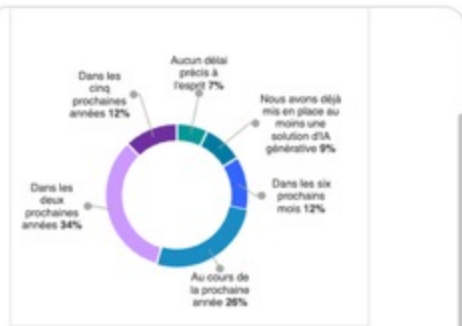
LightOn-Document\_d-Enr...

Attach existing file(s) Upload new file(s)

Company workspace > LightOn-Document\_d-En...  

Citations **1** 2 3 4 5 - 100% + < 21 / 219 >



Notes : la somme n'est pas égale à 100 en raison des arrondis  
Sources : Generative AI survey, mai 2023

**Perspectives d'évolution de l'IA générative au sein de l'entreprise**  
: KPMG, Generative AI: From buzz to business value, An exclusive KPMG survey shows how top leaders are approaching this transformative technology (traduction libre)

adance à l'accélération de l'adoption de l'IA générative par les acteurs que se traduit dans les estimations de taille du marché même de l'IA e. Celui-ci devrait générer un chiffre d'affaires de 1,3 milliard de milliards de :représenter 12 % de toutes les dépenses technologiques d'ici 2032<sup>44</sup>.

42

# More complex workflows



You asked (using OCR)

make a graph illustrating the total vehicle production capacity for each site

| Current Installed Annual Vehicle Capacity |                   |          |                  |
|---|-------------------|----------|------------------|
| Region                                    | Model             | Capacity | Status           |
| California                                | Model S / Model X | 100,000  | Production       |
|   | Model 3 / Model Y | >550,000 | Production       |
| Shanghai                                  | Model 3 / Model Y | >950,000 | Production       |
| Berlin                                    | Model Y           | >375,000 | Production       |
| Texas                                     | Model Y           | >250,000 | Production       |
|   | Cybertruck        | >125,000 | Production       |
| Nevada                                    | Tesla Semi        | -        | Pilot production |
| Various                                   | Next Gen Platform | -        | In development   |
| TBO                                       | Roadster          | -        | In development   |

7BIXJ.png

75.0 KB



OCR executed

[Details](#)



Agent answered

Based on the text extracted from the image using OCR, I can provide a graph illustrating the total production capacity of vehicles for each site.



execute\_python executed Success

[Details](#)

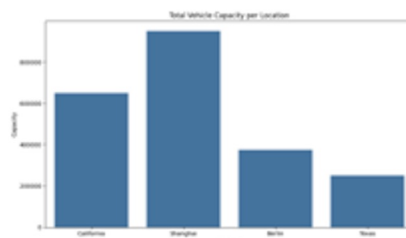


chart.png

21.0 KB

# This workflow calls on a range of tools

Optical Character Recognition



You asked (using OCR)

make a graph illustrating the total vehicle production capacity for each site

| Region     | Model             | Capacity | Status           |
|------------|-------------------|----------|------------------|
| California | Model S / Model X | 100,000  | Production       |
|            | Model 3 / Model Y | >550,000 | Production       |
| Shanghai   | Model 3 / Model Y | >950,000 | Production       |
| Berlin     | Model Y           | >375,000 | Production       |
| Texas      | Model Y           | >250,000 | Production       |
|            | Cybertruck        | >125,000 | Production       |
| Nevada     | Tesla Semi        | -        | Pilot production |
| Various    | Next Gen Platform | -        | In development   |
| TBD        | Roadster          | -        | In development   |

7BIXJ.png  
75.0 KB



OCR executed

Details



Agent answered

Based on the text extracted from the image using OCR, I can provide a graph illustrating the total production capacity of vehicles for each site.



execute\_python executed

Success

Details

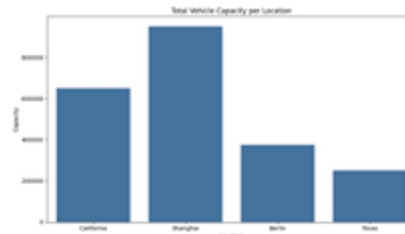


chart.png  
21.0 KB

Python execution

One more thing...



# One more thing...

## Can you use optics to train a (very) large digital model ?

Wang, Muller, et al., arXiv:2409.12965(2024)



# Still plenty of interesting research questions

Platform

**PARADIGM**



*The LLM platform  
for the Enterprise*

Software



*Large Language Models*

Hardware as internal R&D



System

AI software

+

Photonic hardware



Component

Photonic hardware



Lab  
experiment

Back to  
the lab

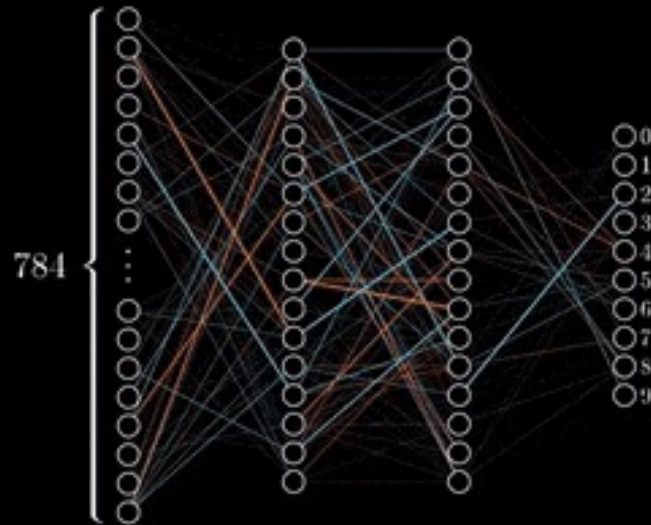


# Training with backpropagation

FORWARD : inference



Training in progress...



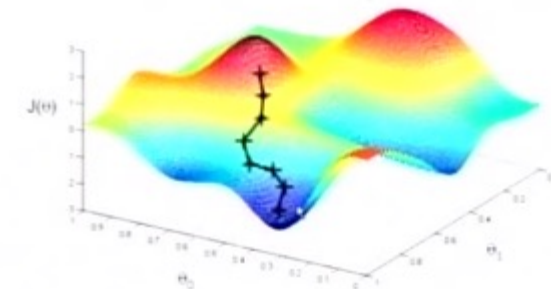
## « **backprop'** »

- de-facto standard for training
- Power-intensive
- Constraints to a layered NN architecture

BACKWARD : Error propagation and weight-tuning based on local gradients

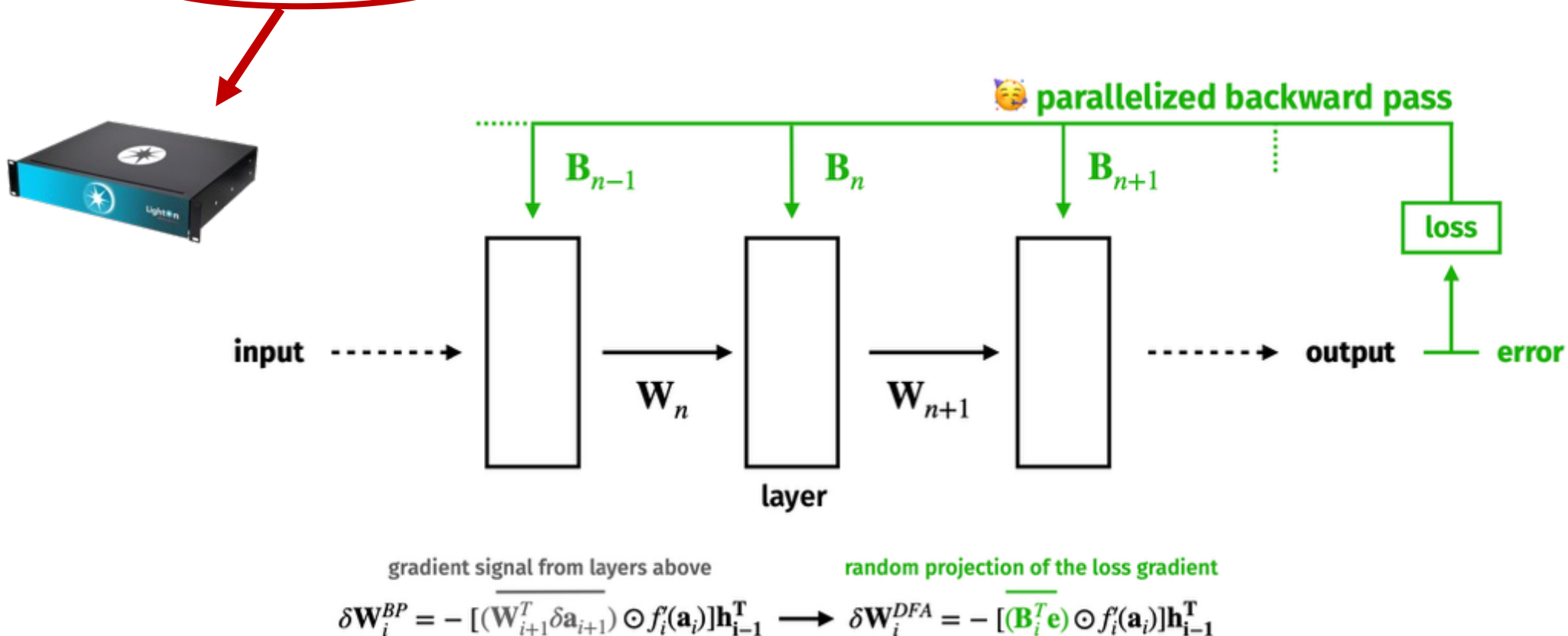


Gradient Descent



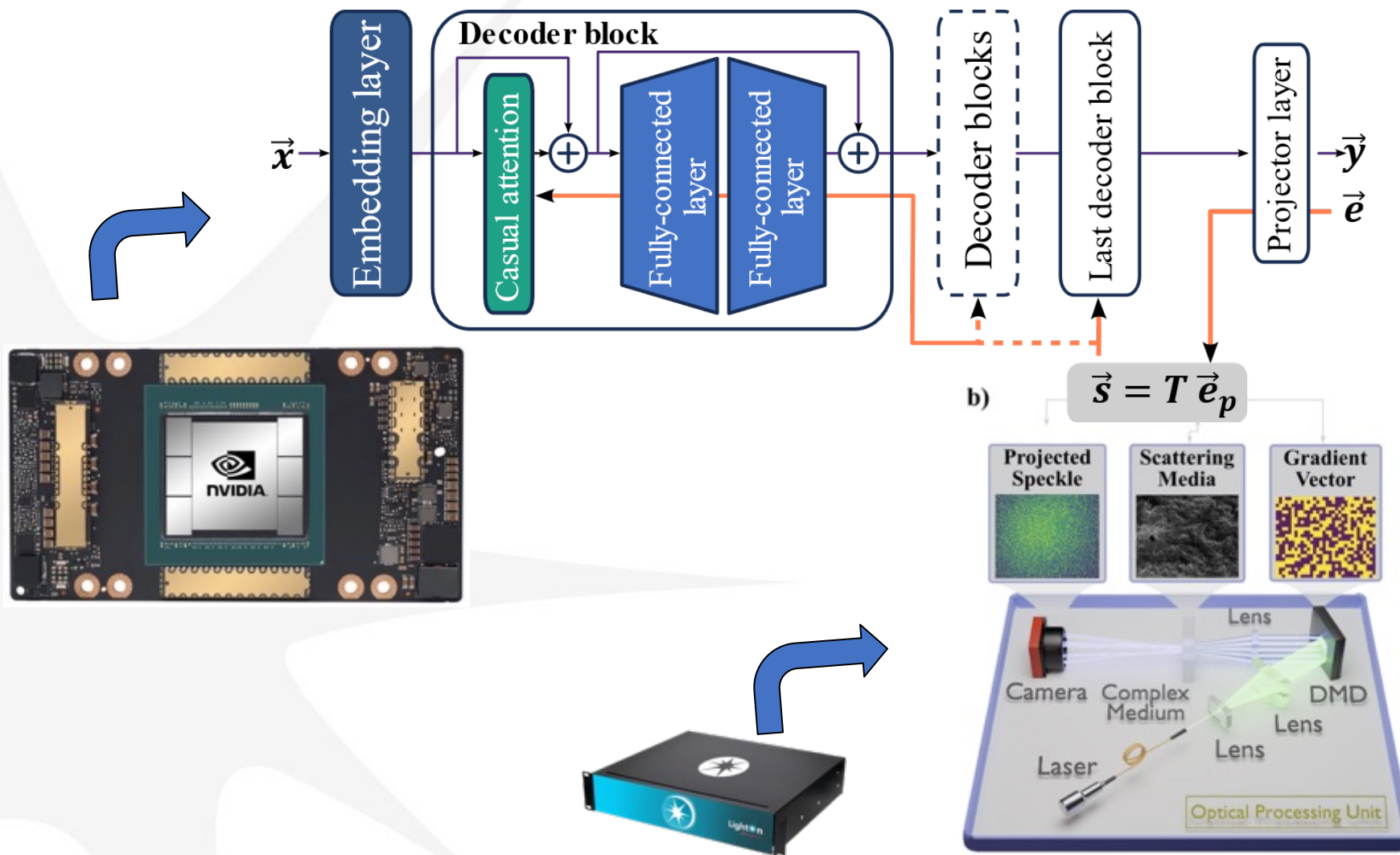
# Training with Direct Feedback Alignment (DFA)

Use **random weights** to **directly** deliver feedbacks from the global loss.



General-purpose: agnostic to model architecture & scale well

# ODFA on Text Transformer



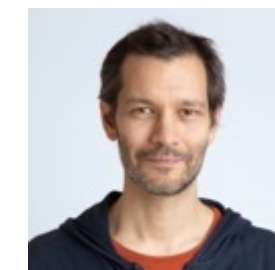
- 1B parameters GPT-like transformer
- 400M parameters directly receive optical signal

Cornell Movie-Dialogs Corpus  
**Example in the Dataset**

- MILES: Back at you.
- JACK: Love you, man.
- MILES: Yeah.
- JACK: So I'll see you at the rehearsal.

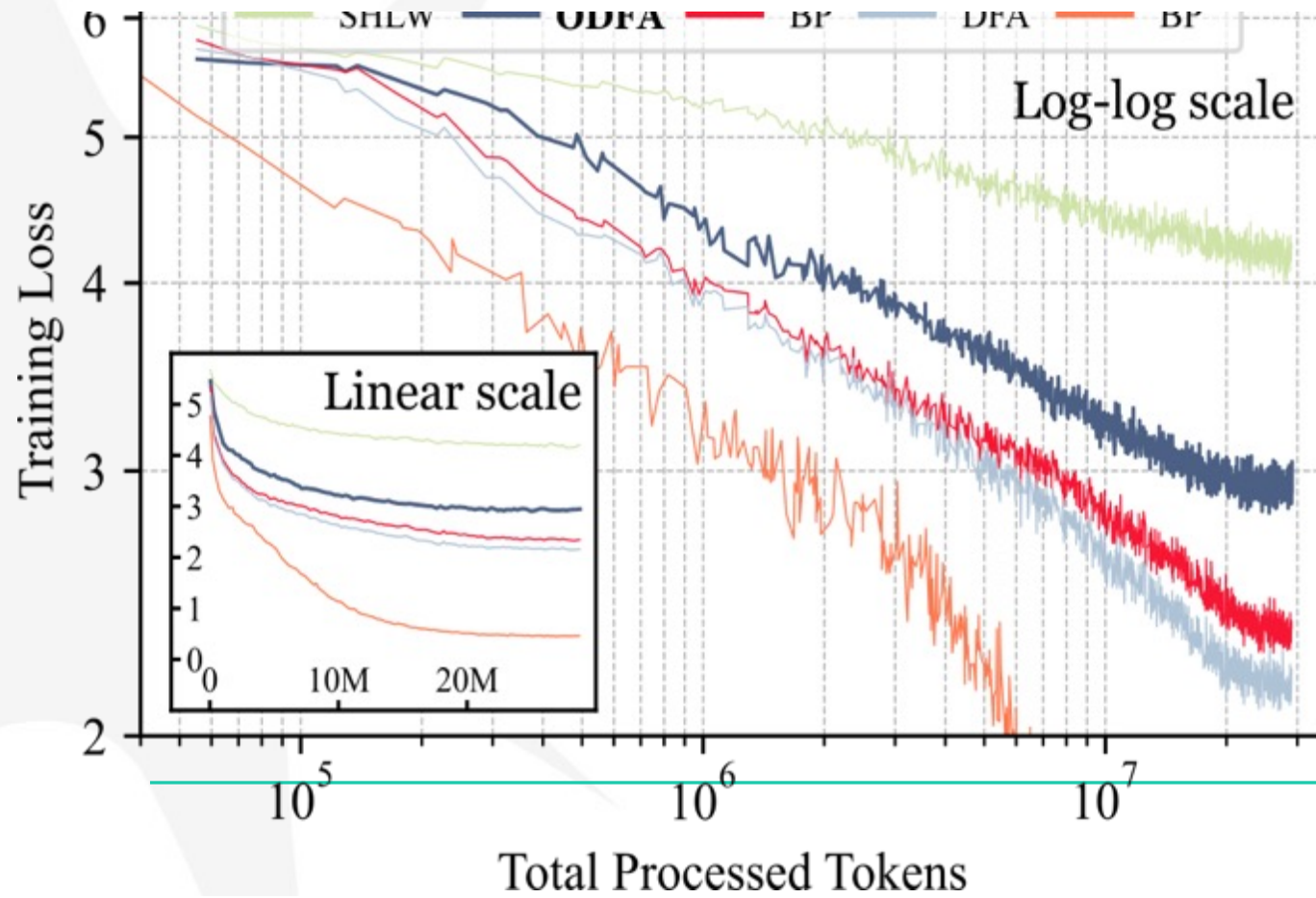


Ziao Wang  
(LKB)



Kilian Müller  
(LightOn)

# ODFA on Text Transformer: Performance



# ODFA on Text Transformer: Performance

## 0% Training

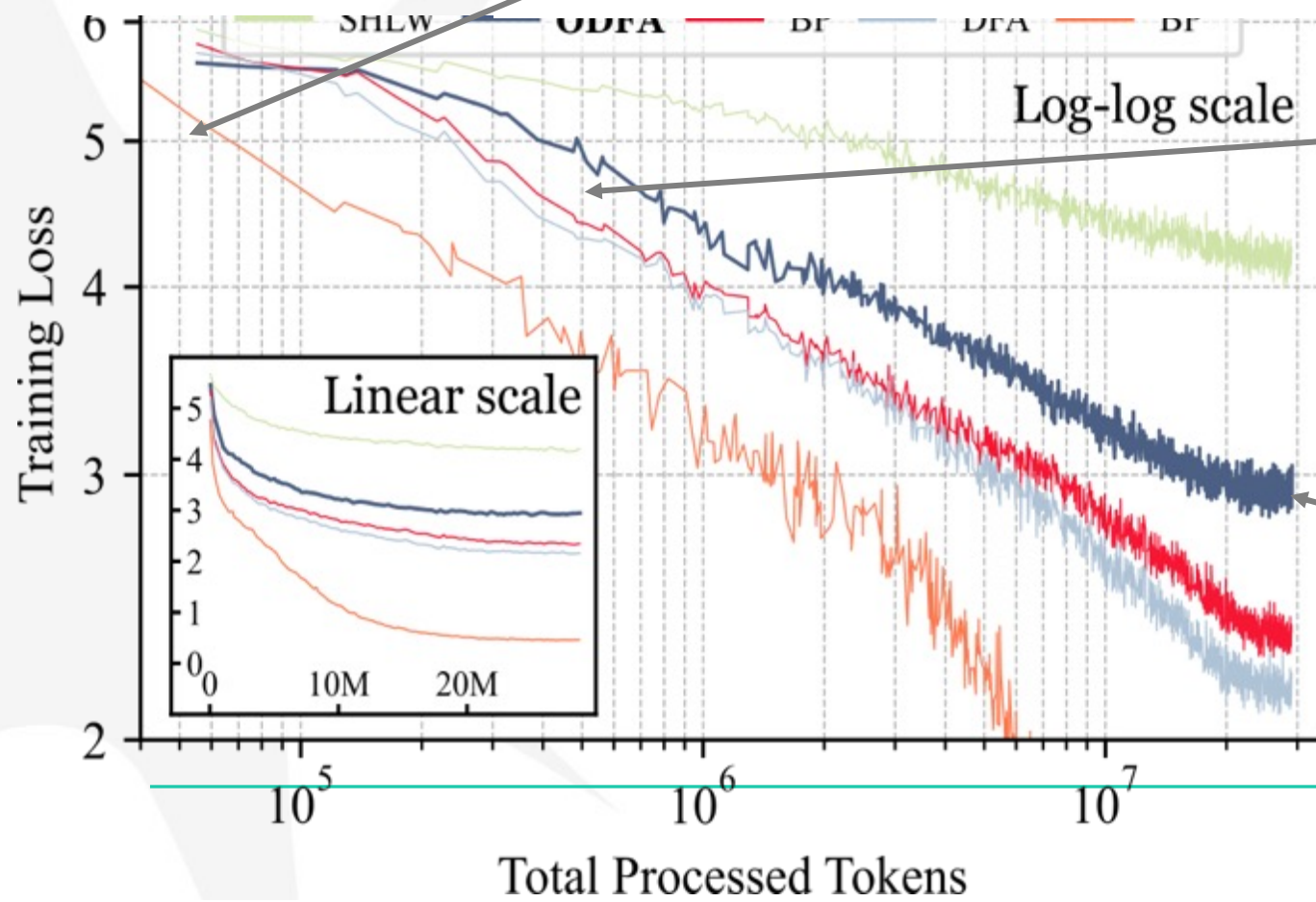
aQ o do wantQain do do twoain o  
re[ through show" happenQ start. 2lt  
JULIANNE JULIANNE

## 20% Training

JACKIE : Snkups to the brags, but I  
know that. ↵  
CRAYC : It ' s just a concer. ↵

## 100% Training

JOE : So what are you talking about? ↵  
ED : I mean I'm sorry, you know how  
to know you are. ↵



# ODFA on Text Transformer: Performance

## 0% Training

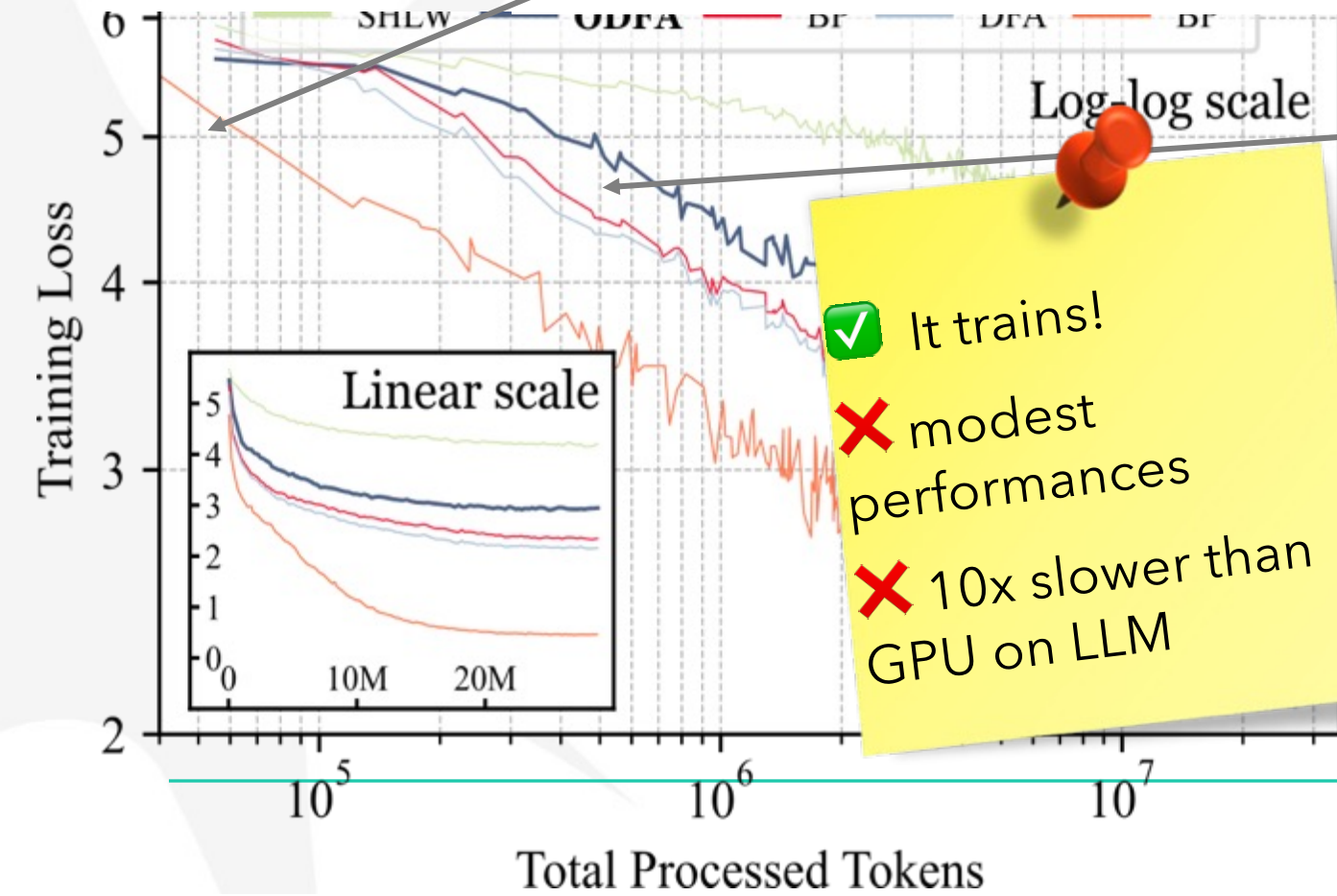
aQ o do wantQain do do twoain o  
re[ through show" happenQ start. 2lt  
JULIANNE JULIANNE

## 20% Training

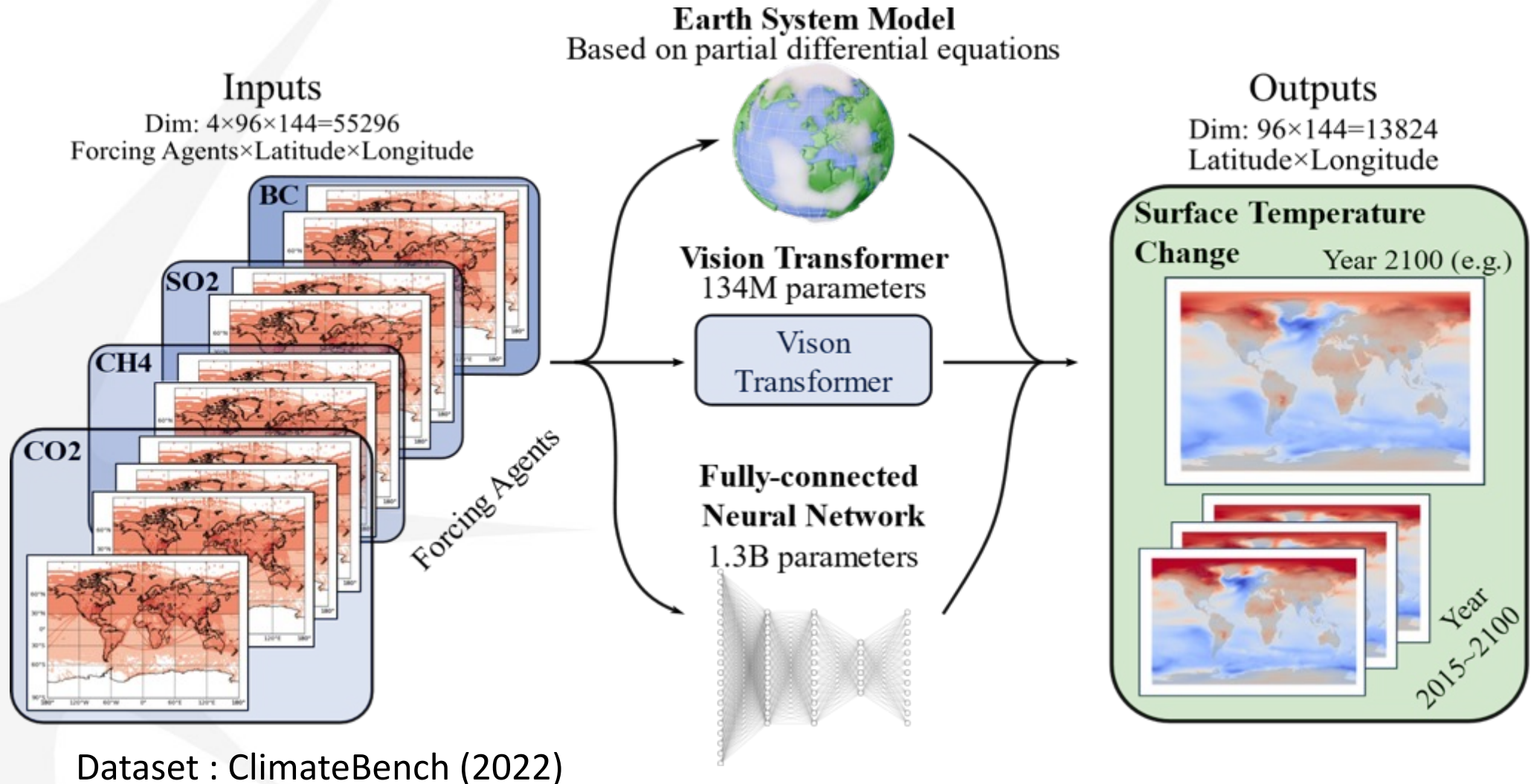
JACKIE : Snkups to the brags, but I  
know that. ↵  
CRAYC : It ' s just a concer. ↵

## 100% Training

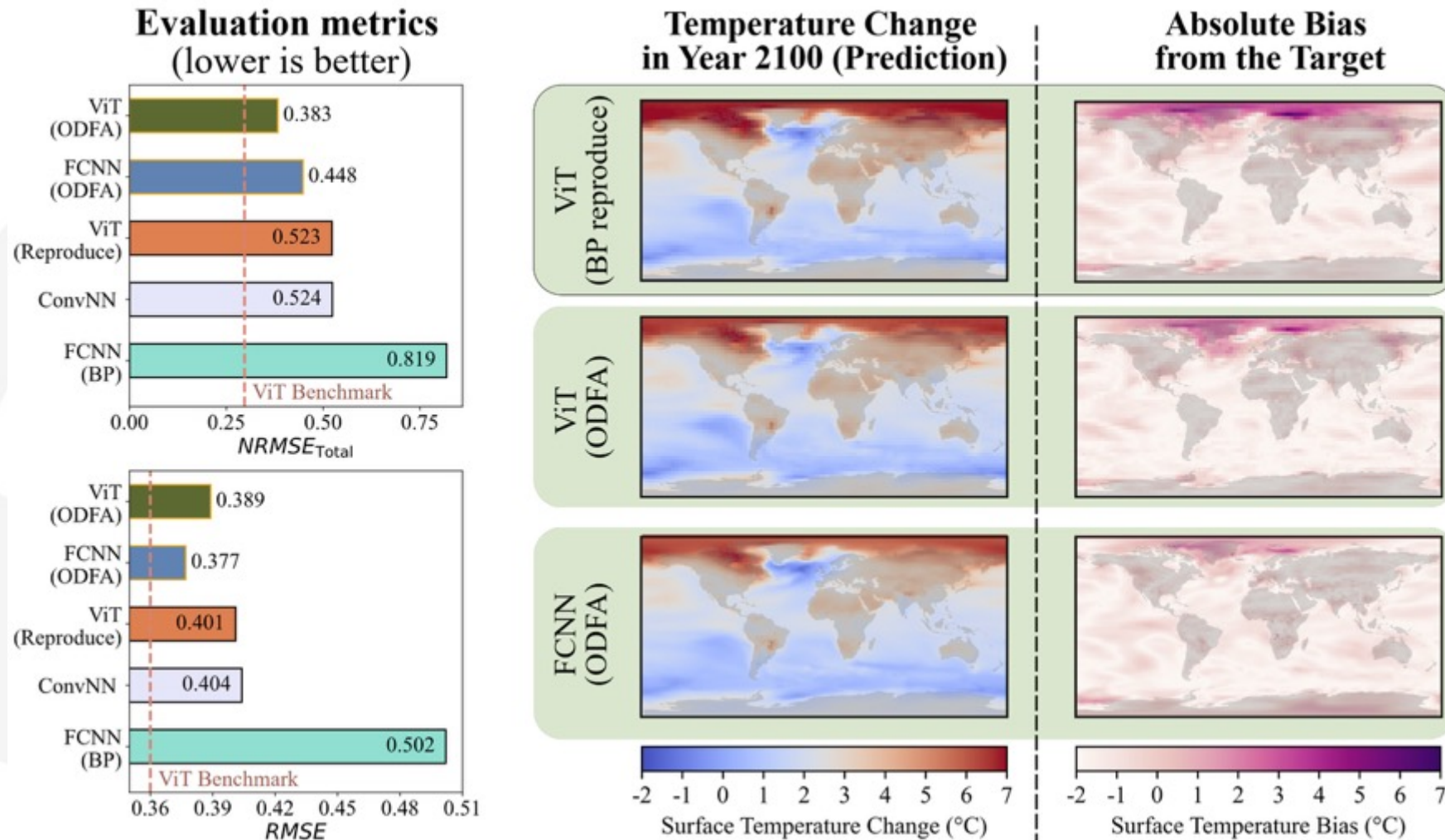
JOE : So what are you talking about? ↵  
ED : I mean I'm sorry, you know how  
to know you are. ↵



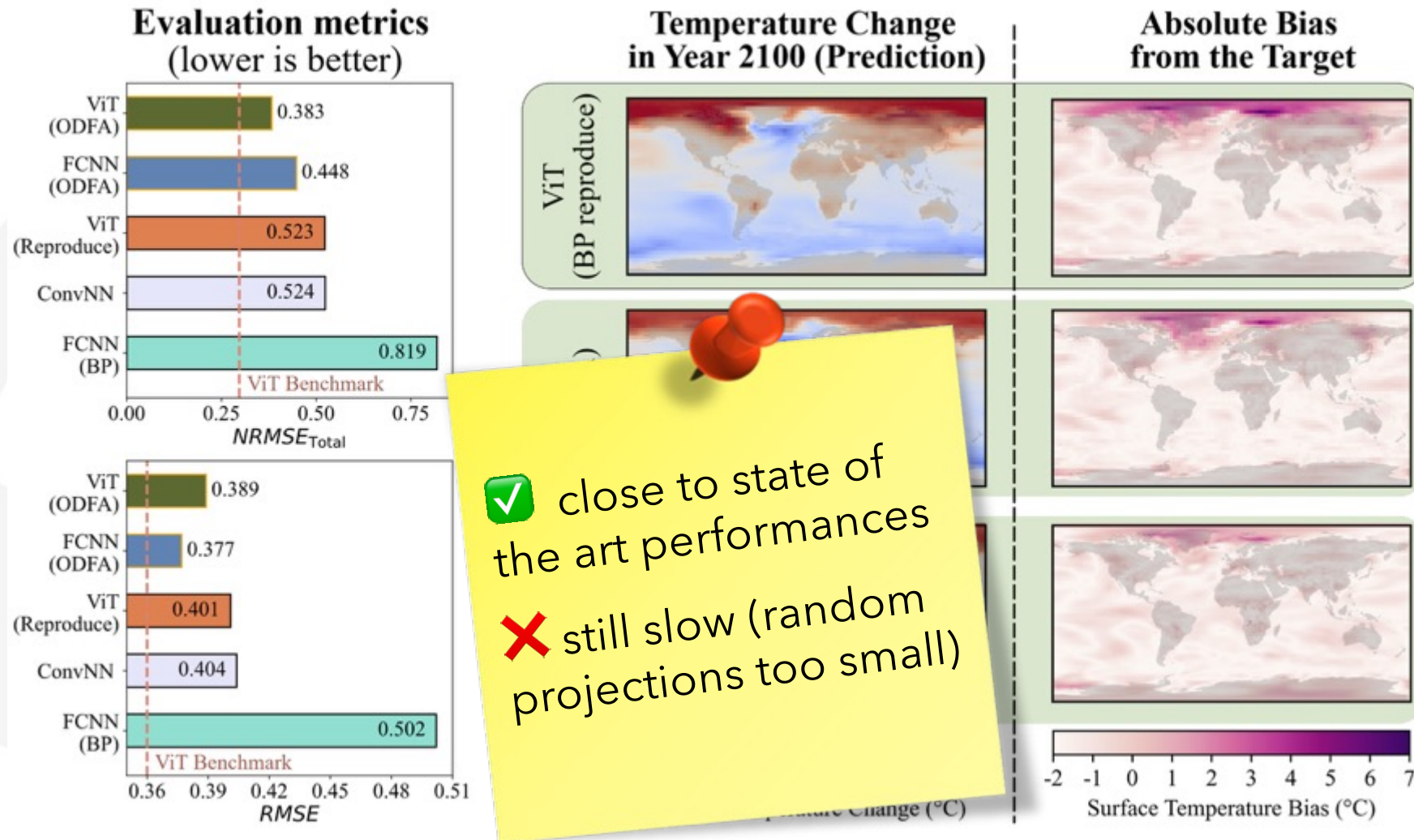
# Vision Task: More Suitable for ODFA



# FCNN and ViT with ODFA on Climate Projection

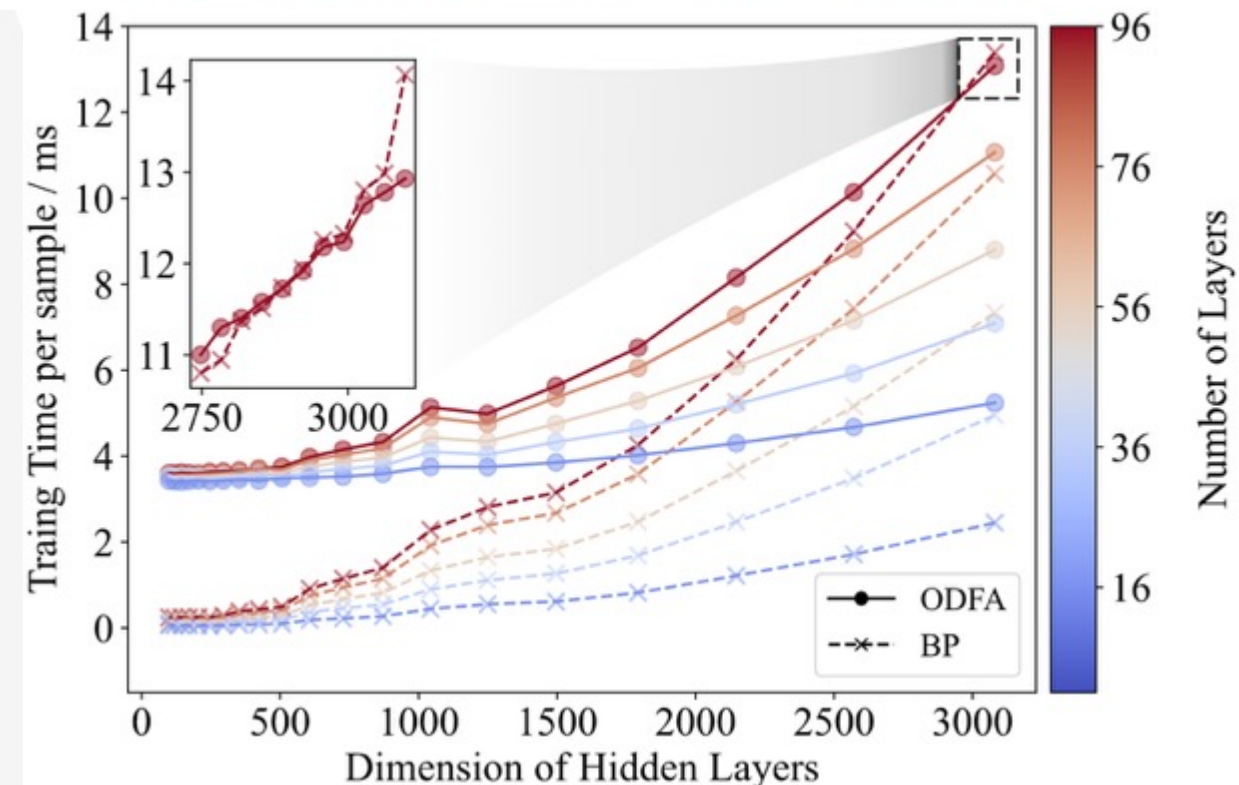


# FCNN and ViT with ODFA on Climate Projection

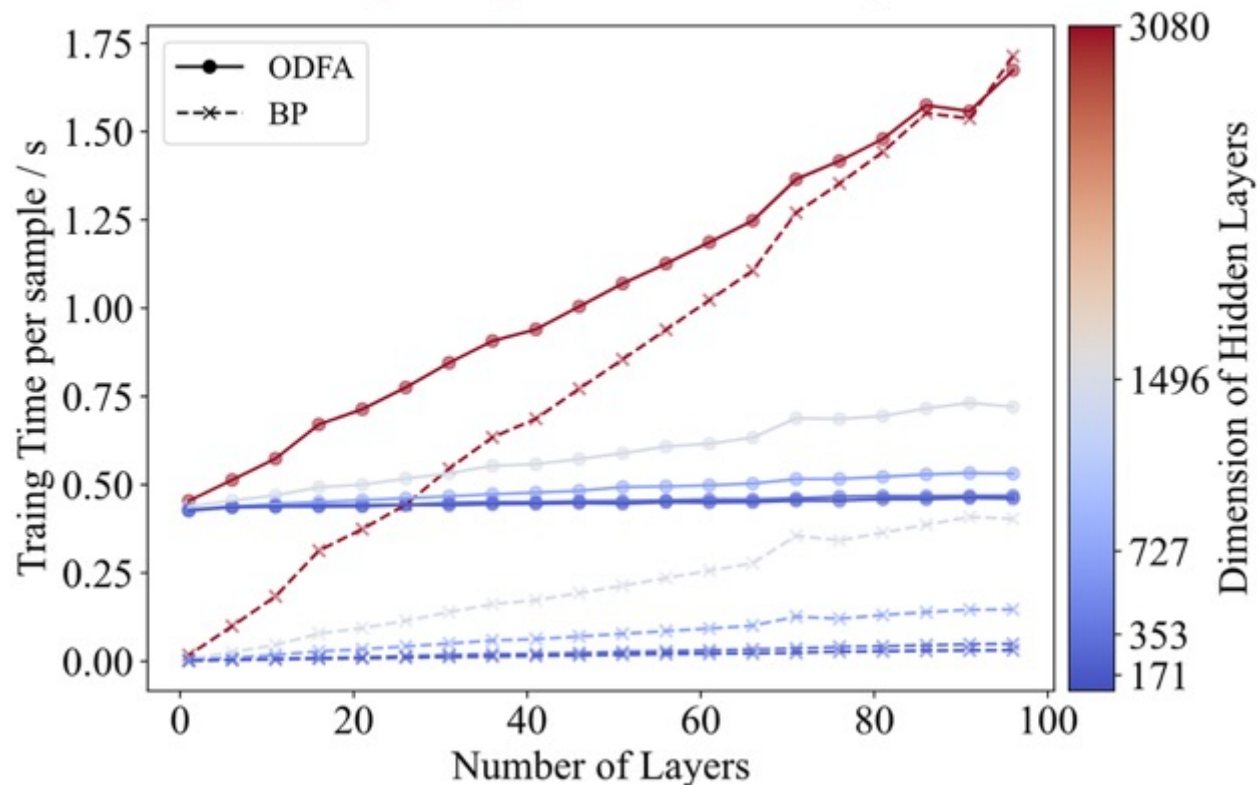


# Towards Extreme Scale: When will ODFA be faster than BP?

Scaling along the dimension of hidden layers

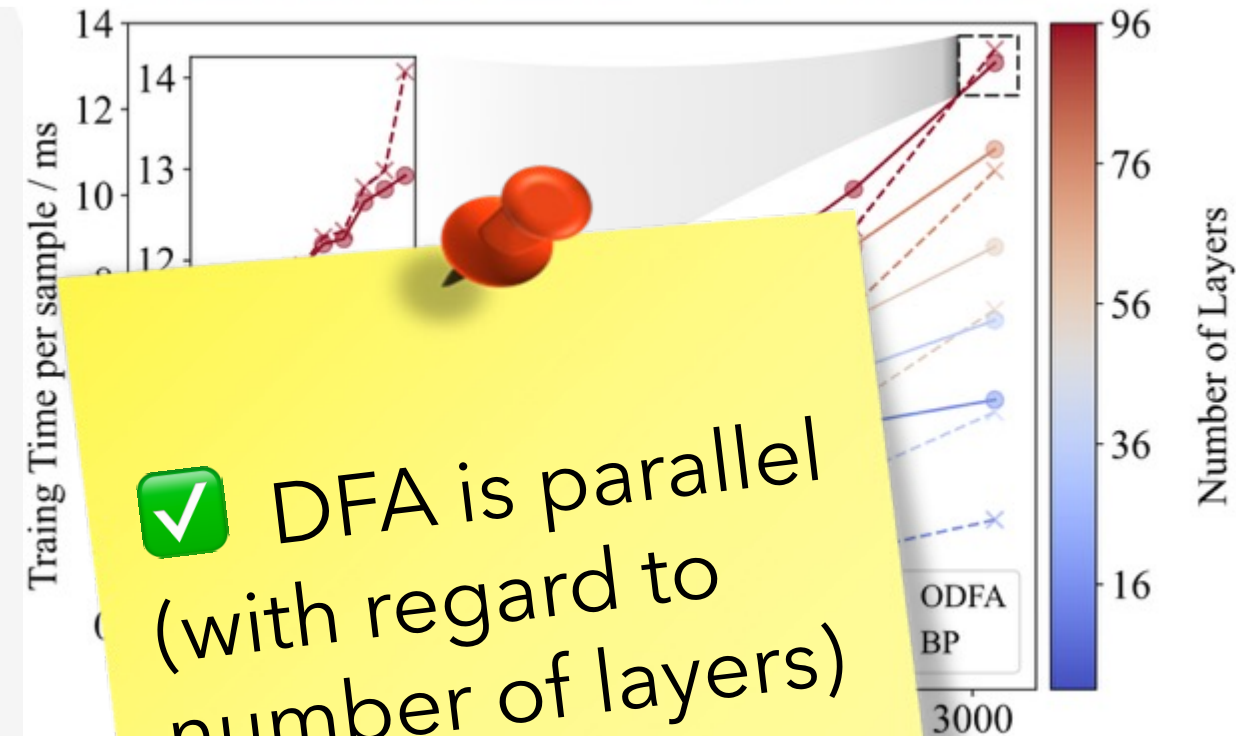


Scaling along the number of layers

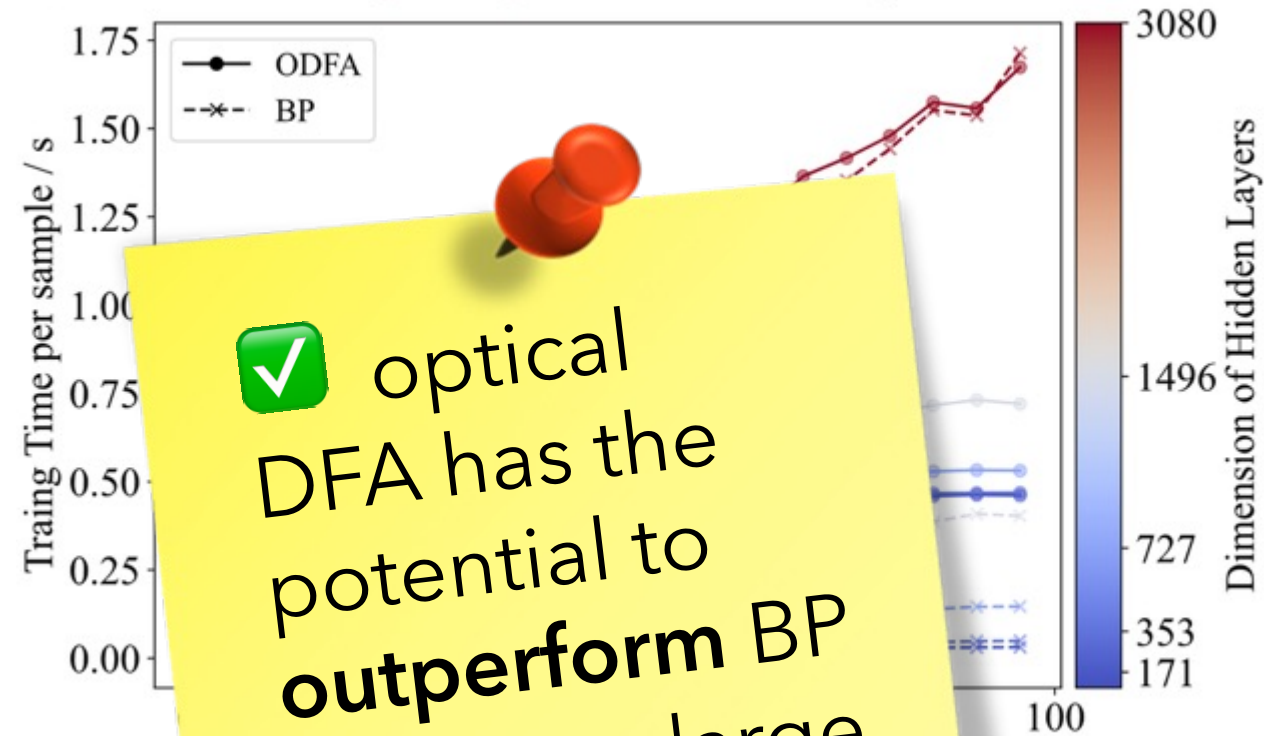


# Towards Extreme Scale: When will ODFA be faster than BP?

Scaling along the dimension of hidden layers



Scaling along the number of layers



By way of conclusion

The first publicly traded GenAI company in Europe !



Nov 26th, 2024

# The entrepreneurial journey

