

# Optical random features for large-scale machine learning

---

**Laurent Daudet**

Paris Diderot University

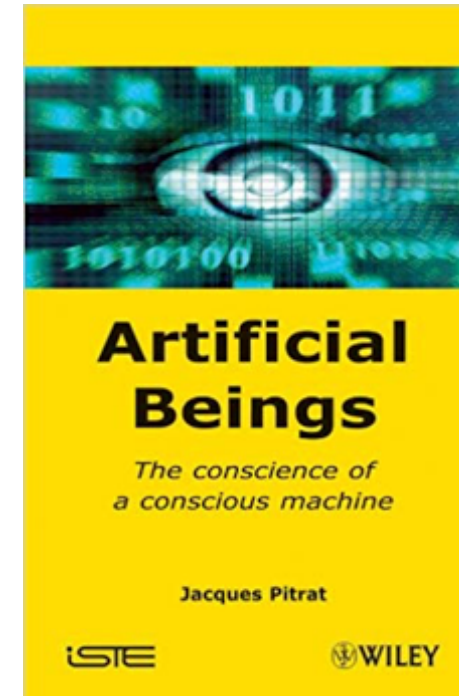
Co-founder & CTO at LightOn

laurent@lighton.ai



# In memoriam Jacques Pitrat (1934-2019)

---



« We must not blindly imitate human behavior because computers are tools which work differently from our brain. »

# The team

---



**Sylvain Gigan**

LKB (UPMC / ENS / CdF)



**Florent Krzakala**

LPS (UPMC / ENS)



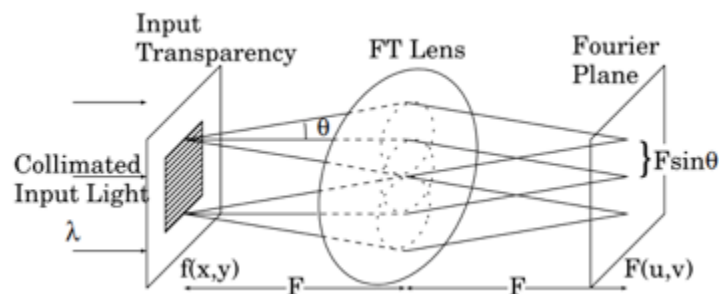
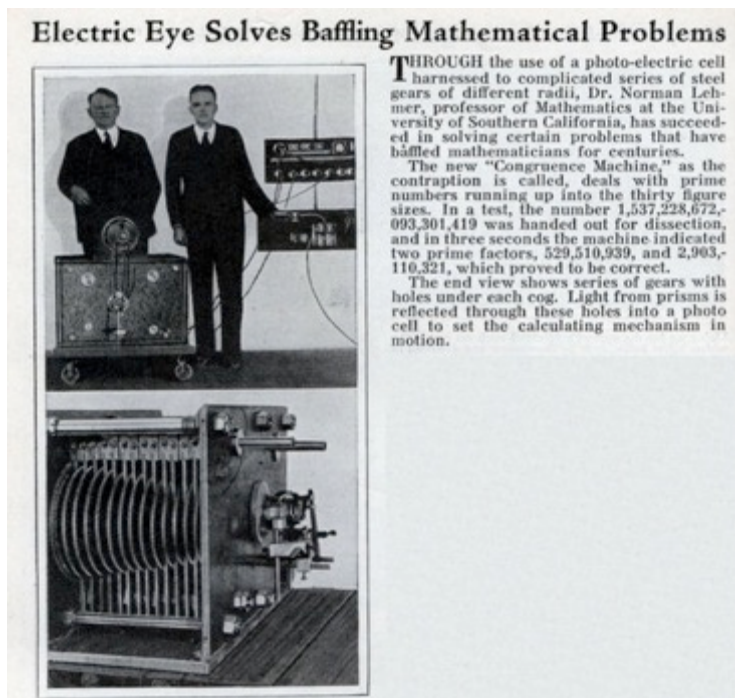
**Igor Carron**

Nuit Blanche / LightOn

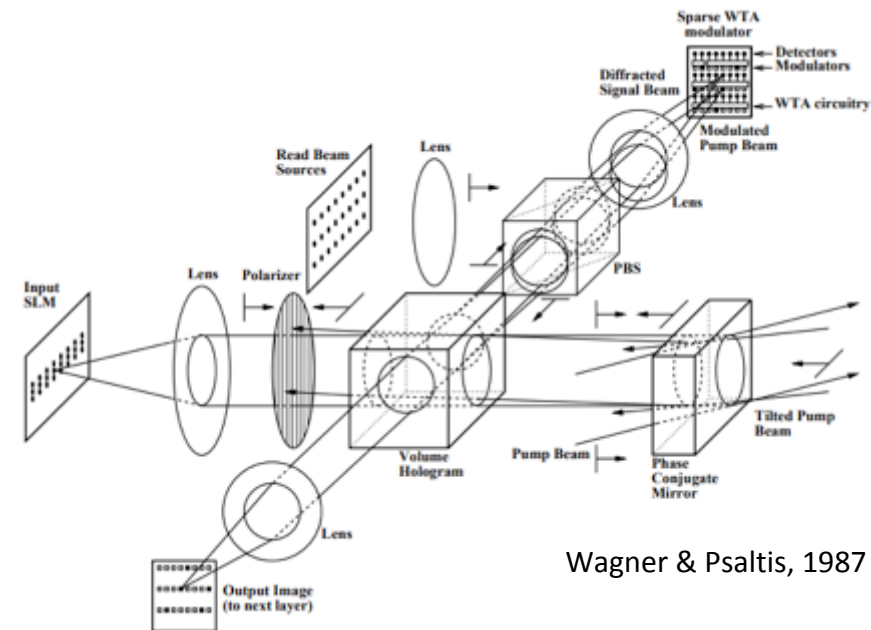
And *many* others in these labs and at LightOn

# A short history of Optical Processing of Information

From Sieves ... to Fourier Transforms ... all the way to Neural Networks



$$F(u,v) = \iint f(x,y) e^{i\frac{2\pi}{\lambda F}(xx' + yy')} dx dy$$



Wagner & Psaltis, 1987

1930's

1950's

1980's



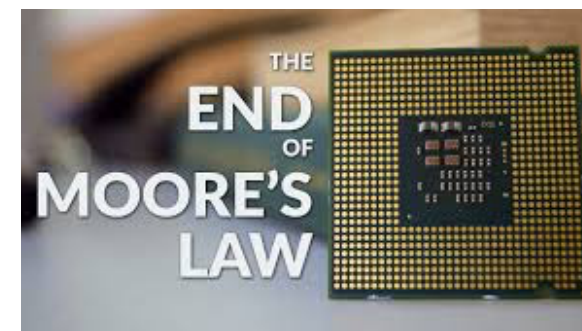
Then  
Came  
Winter



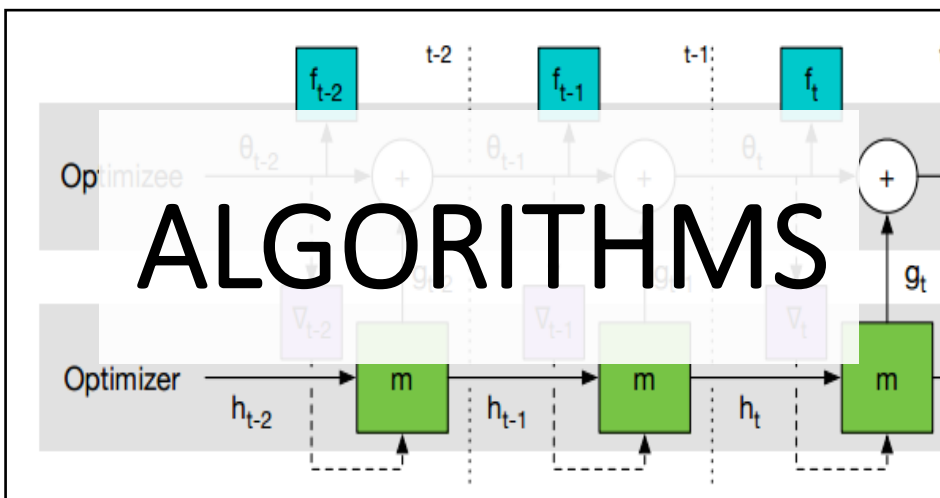
# Rebooting Optical Computing: the AI era



SCALABLE ?



<https://www.youtube.com/watch?v=Ak7HPuuJ1Ow>



SUSTAINABLE ?



11 Dec 2017

Information  
theory  
Compressive  
Sensing



**Igor**

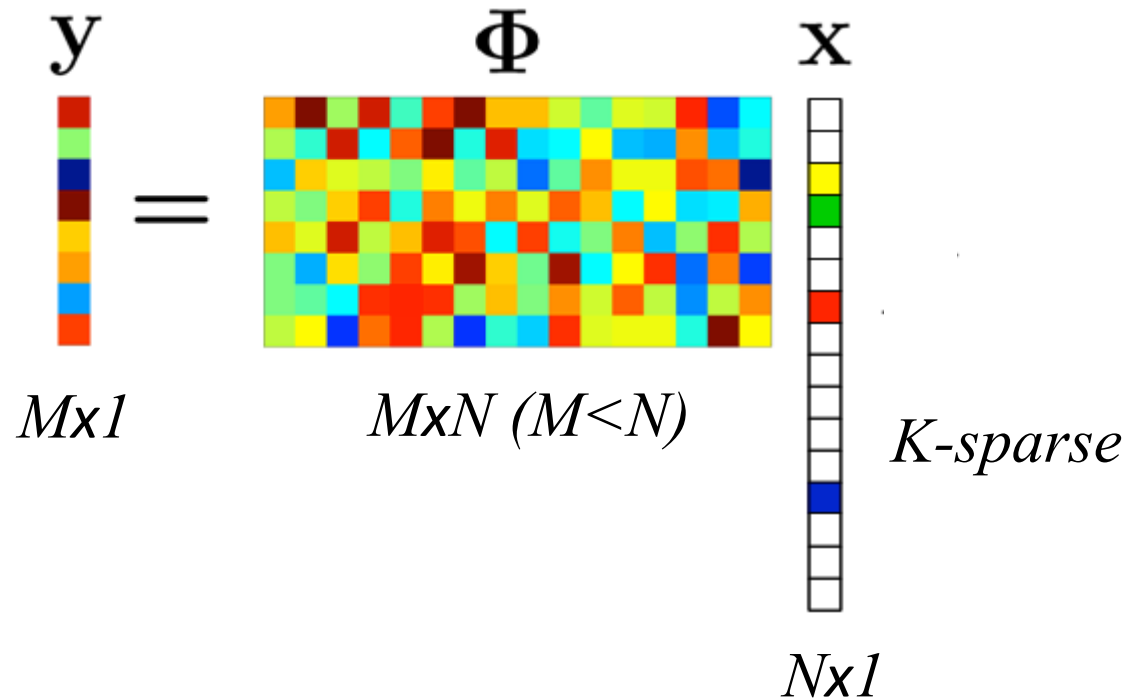


<http://nuit-blanche.blogspot.com>



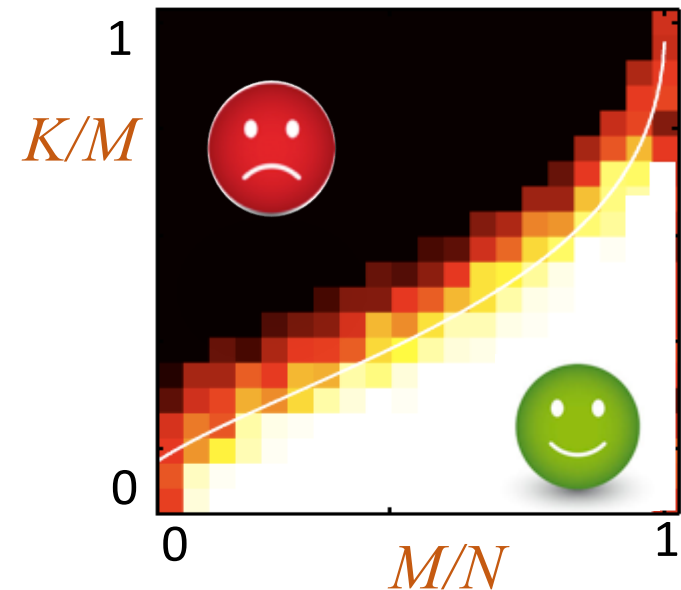
**Laurent**

# Compressive Sensing



Can one recover  $x$  from  $y$  ?

YES with tractable algorithms  
for right values of  $N$ ,  $M$ ,  $K$



# Lessons from Compressive Sensing

---

- Signals can be sampled at the level of their information content
- Random Projections are very good for sensing at low data rate
- Strong theoretical background and large empirical evidence



Information  
theory  
Compressive  
Sensing



**Igor**



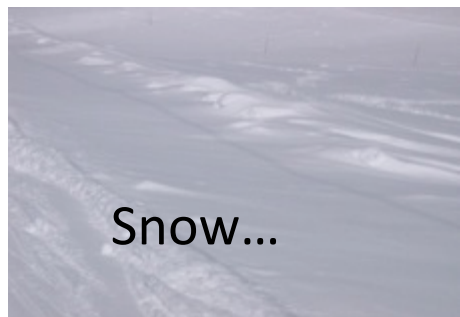
**Laurent**

Light Transport  
in Diffusive  
Media



**Sylvain**

# Light transport in diffusive media

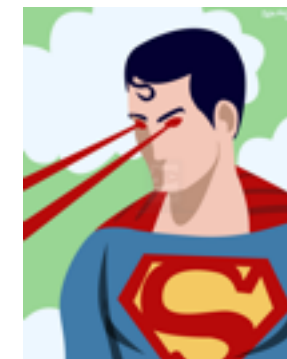


How deep can one see ?

Not much



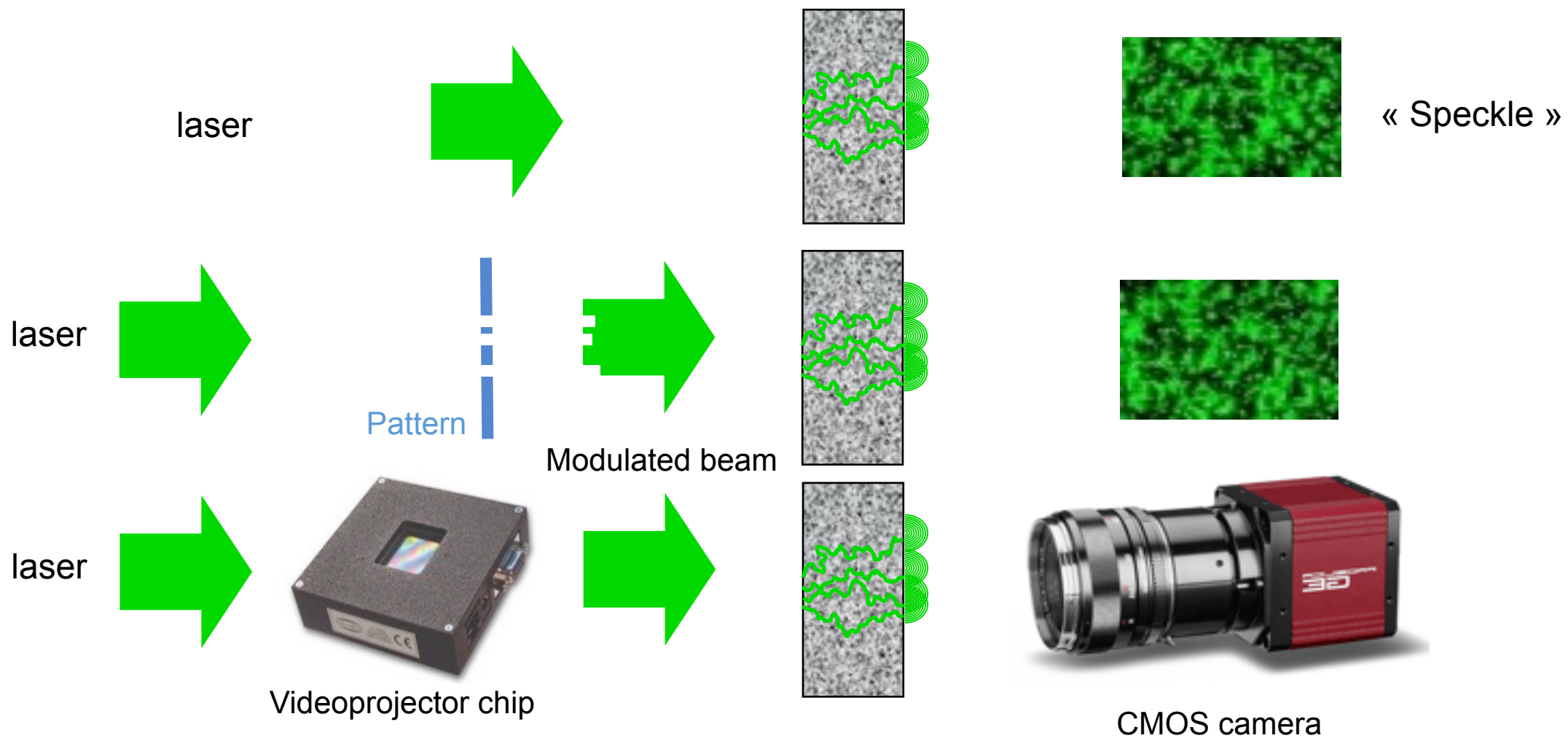
or



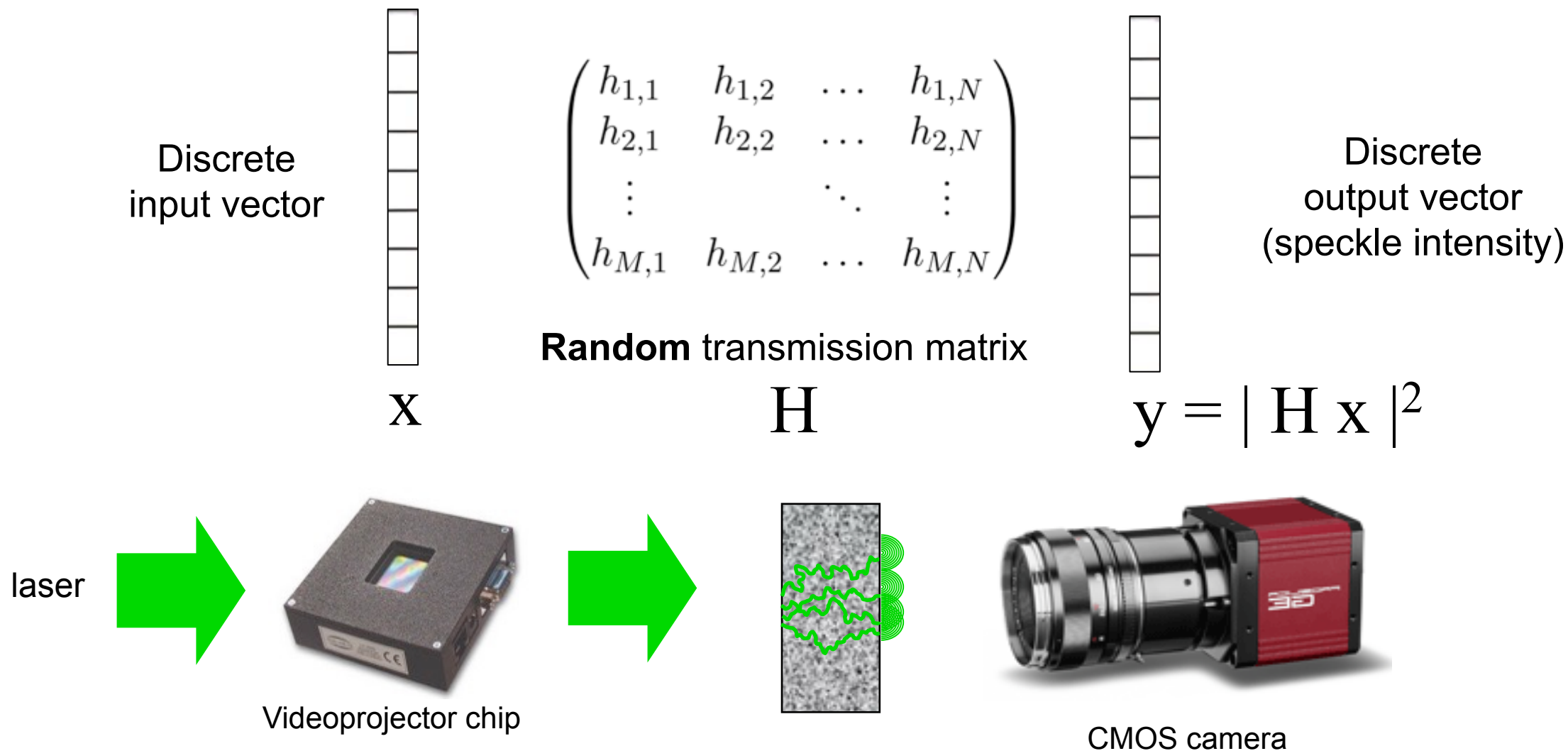
?



# Scattering: a coherent process



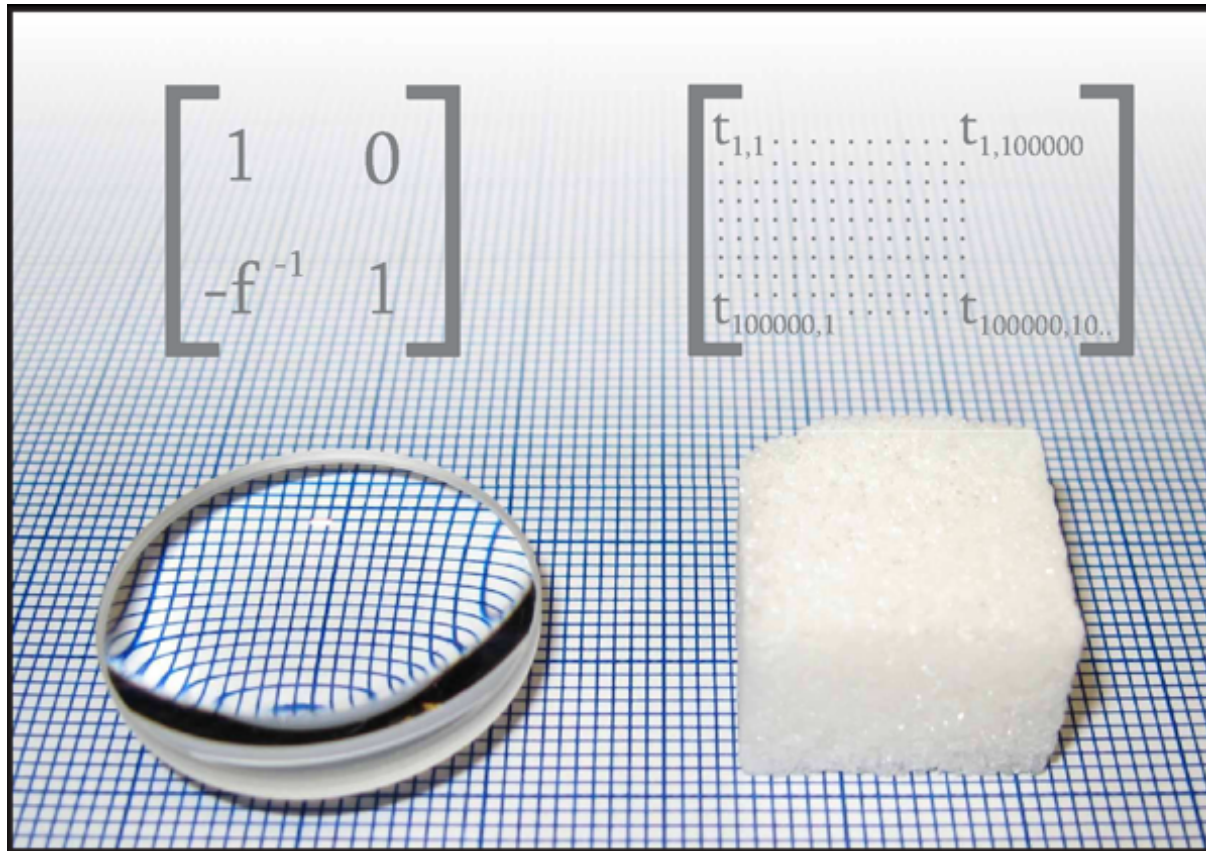
# Scattering: a coherent process



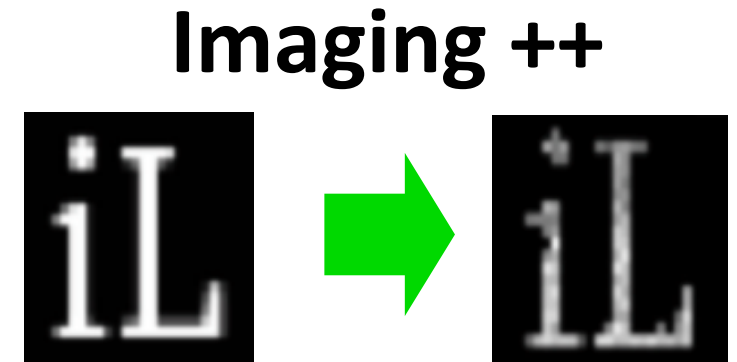
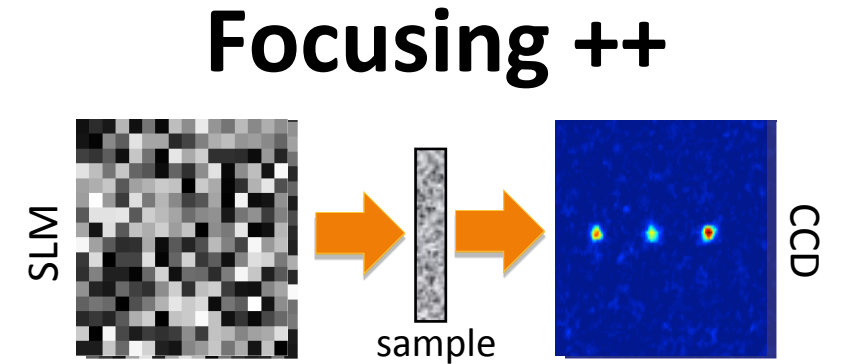


# The transmission matrix

Scattering materials are « super-lenses »



Popoff et al. Phys. Rev. Lett. 104,100601 (2010)



Liutkus et al., Scientific Reports 4, 5552 (2014)



# Lessons from Light Transport in Diffusing Media

---

- Scattering preserves the information content
- Scattering *optimally scrambles* information
  - just like a Random Projection
  - just like in Compressive Sensing
- Matrix-vector multiplication, followed by non-linearity: sounds familiar ?

Information  
theory  
Compressive  
Sensing



**Igor**

Machine  
Learning

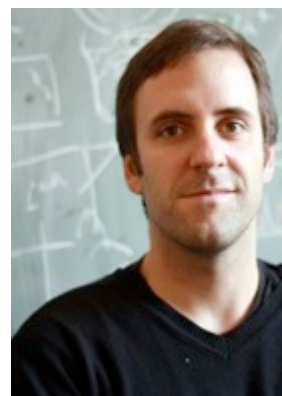


**Florent**



**Laurent**

Light Transport  
in Diffusive  
Media



**Sylvain**

# Random Projections in Machine Learning

- 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

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

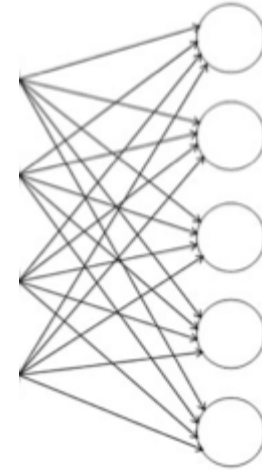
then for any set  $A$  of  $n$  points  $\in \mathbb{R}^d$  there exists a map  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$  such that for all  $x_i, x_j \in A$

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

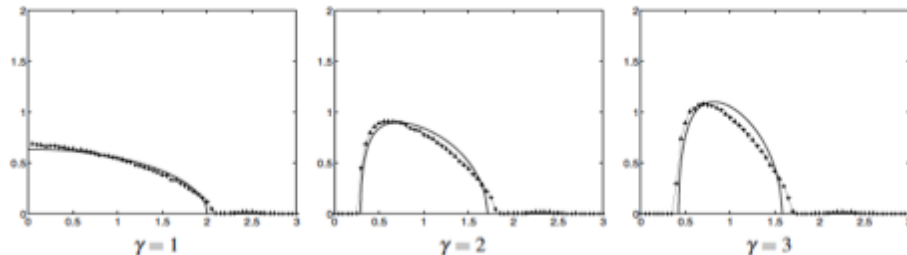


- NIPS 2017 Test of Time Award  
“*Random Features for Large-scale Kernel Machines*”, Rahimi, Recht, 2008

# Random Projections in Machine Learning



- A matrix-vector multiplication followed by a non-linearity: a fully connected layer of a Neural Network
- Fixed dense random weights - you can guarantee their distribution (Gaussian iid complex)



Marčenko-Pastur law on  
singular values

- Random projections made  $O(n^2) \rightarrow O(1)$

# Lessons from Random Projections in Machine Learning

---

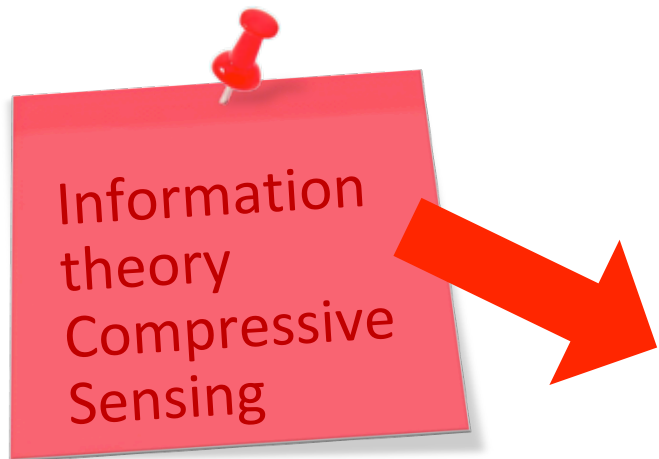
- Random projections act as dimensionality reduction or expansion
- Supervised or unsupervised
- Can also be seen as a fixed dense layer in a Deep Learning model



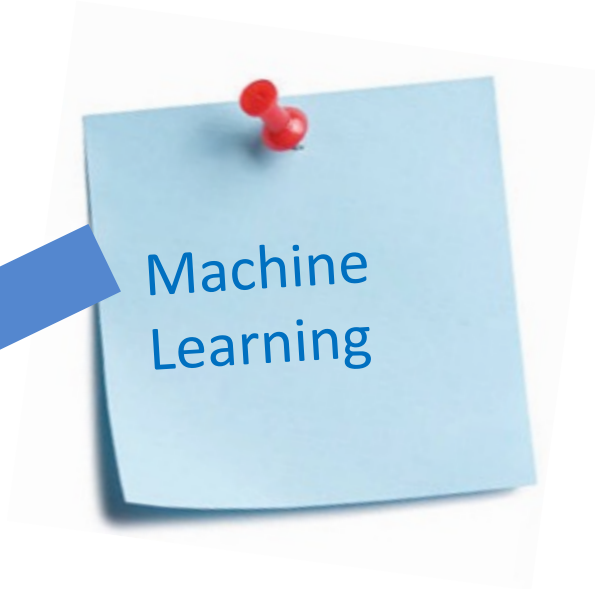
« Ask not what AI can do for Physics  
– ask what Physics can do for AI »

# The Convergence

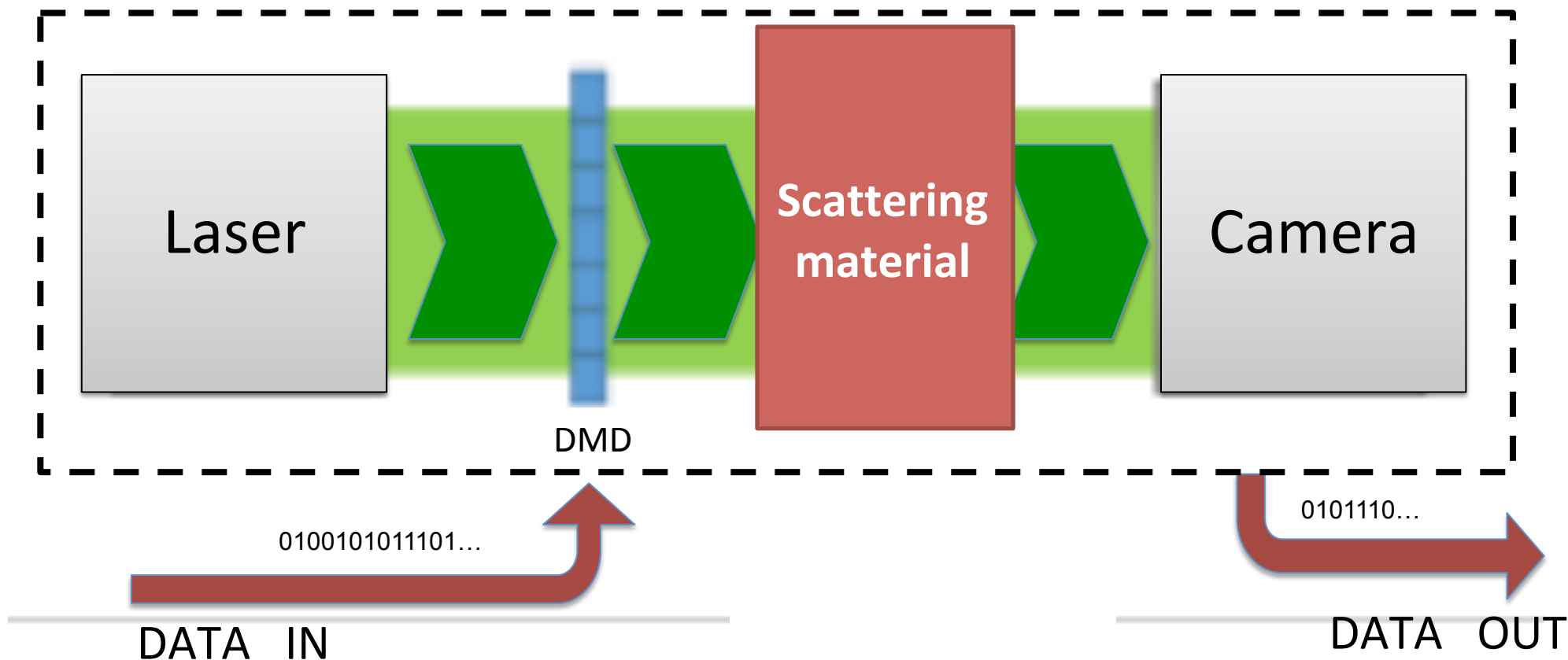
---



**Optical Processing  
for Large Scale ML**



# Optical Processing for Large Scale ML



## Optical Processing Unit (OPU)

# Optical Processing for Large Scale ML

The OPU performs **Random Projections** in the analog domain

$$y = |Hx|^2$$

with H a complex random iid matrix



EXTRA-LARGE

&

SUPER-FAST

H of size higher than  
 $10^6 \times 10^6$   
(TBs of memory)

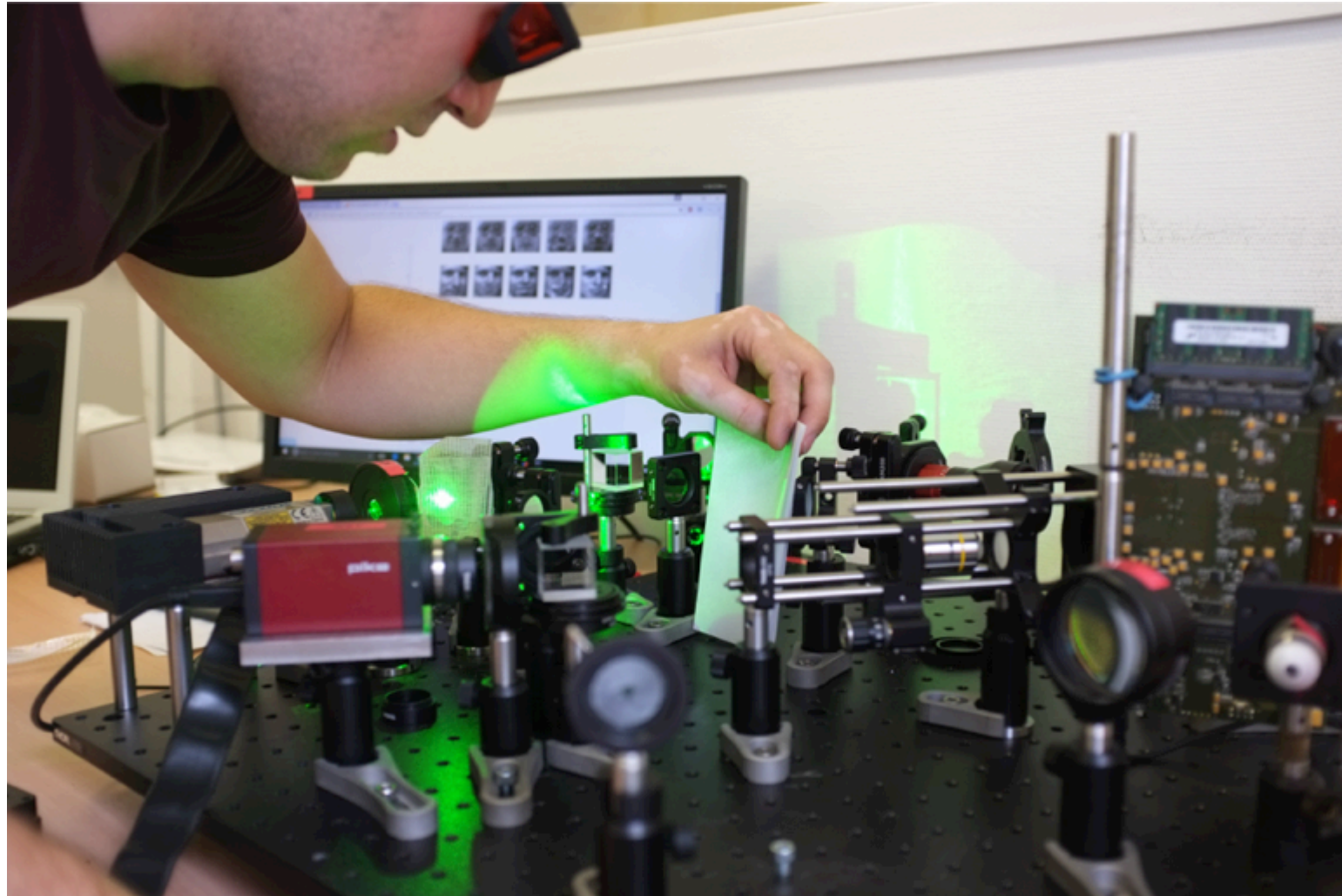
kHz operation  
 $\rightarrow 10^3$  such  
multiplies / s



Equivalent  $10^{15}$  operations\* / s ... for a few W

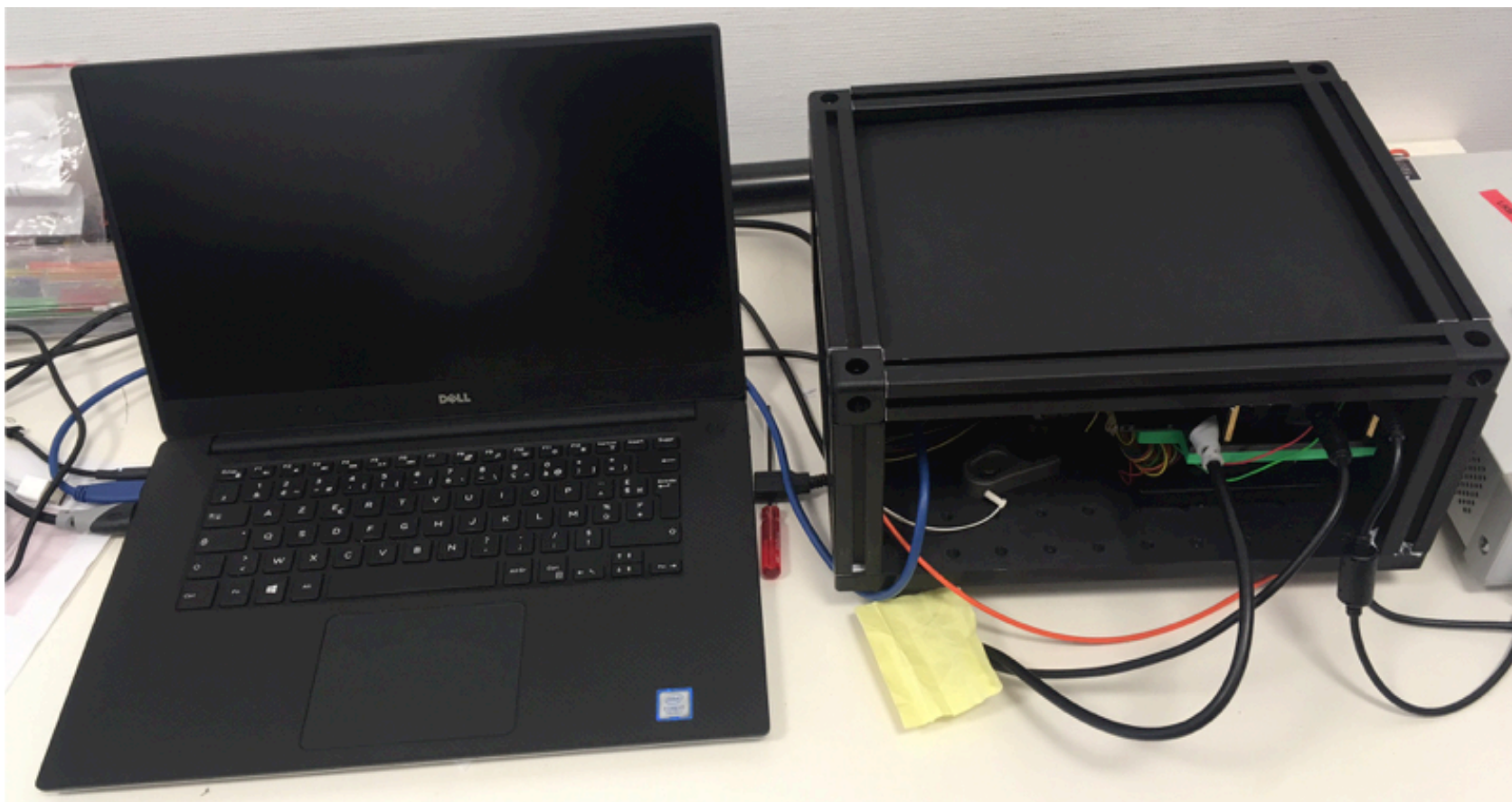
\* Analog « operations » not directly comparable to flops

# Optical Processing for Large Scale ML



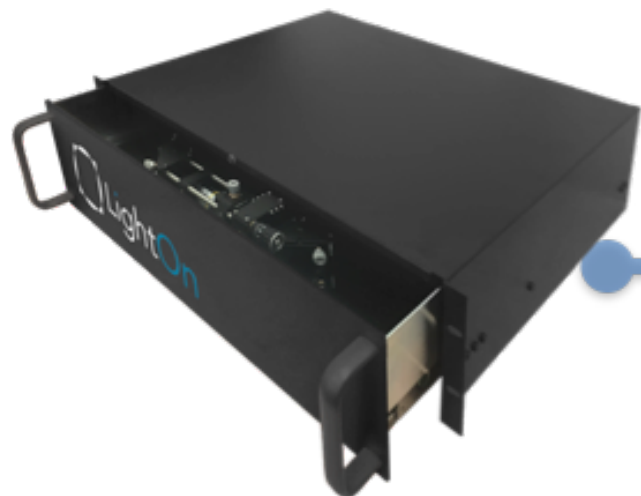


# Optical Processing for Large Scale ML

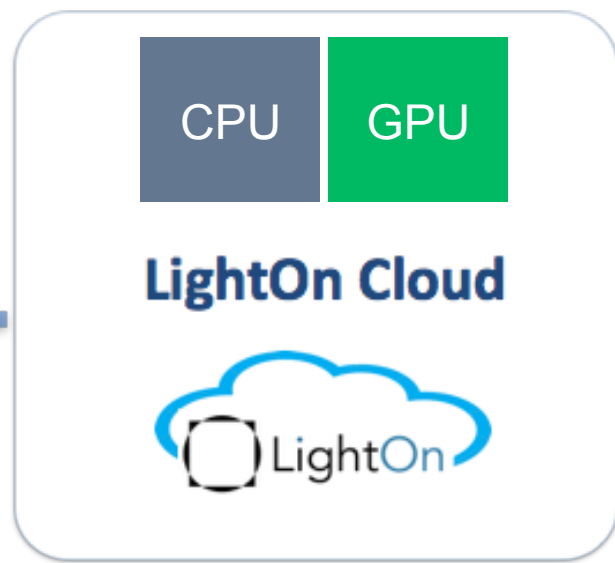


Spring 2017 - The first « OPU »: Optical Processing Unit

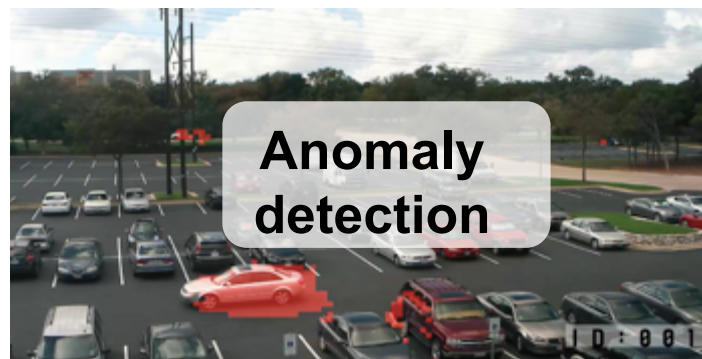
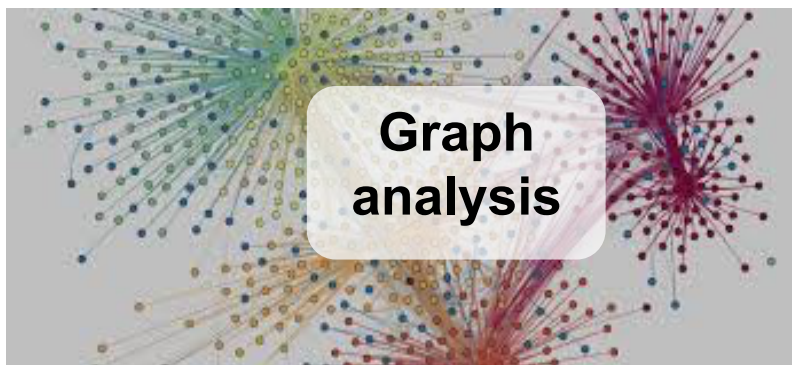
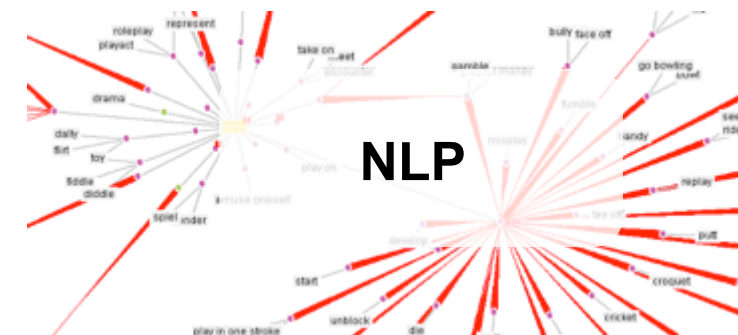
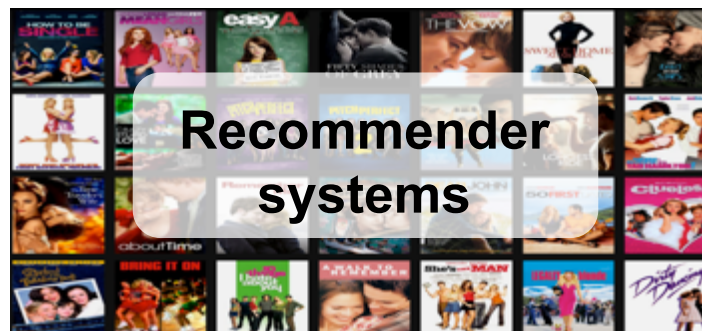
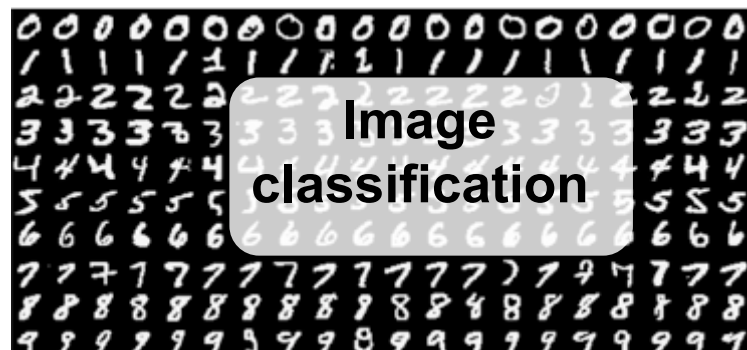
# Optical Processing for Large Scale ML



« Zeus » OPU prototype  
30 W



# Some use cases



# Case study 1: classification with kernel ridge regression

---

training

$$\underset{\beta \in \mathbb{R}^{p \times q}}{\operatorname{argmin}} \quad \overset{\text{U : data}}{\mathbf{U}} \overset{\text{Y: labels}}{\beta} - \mathbf{Y} \quad \|\mathbf{U}\beta - \mathbf{Y}\|_2^2 + \gamma \|\beta\|_2^2$$

Example : classifying the MNIST database

training set of 60000 training pictures  
(28x28) of handwritten digits

test set of 10000 digits

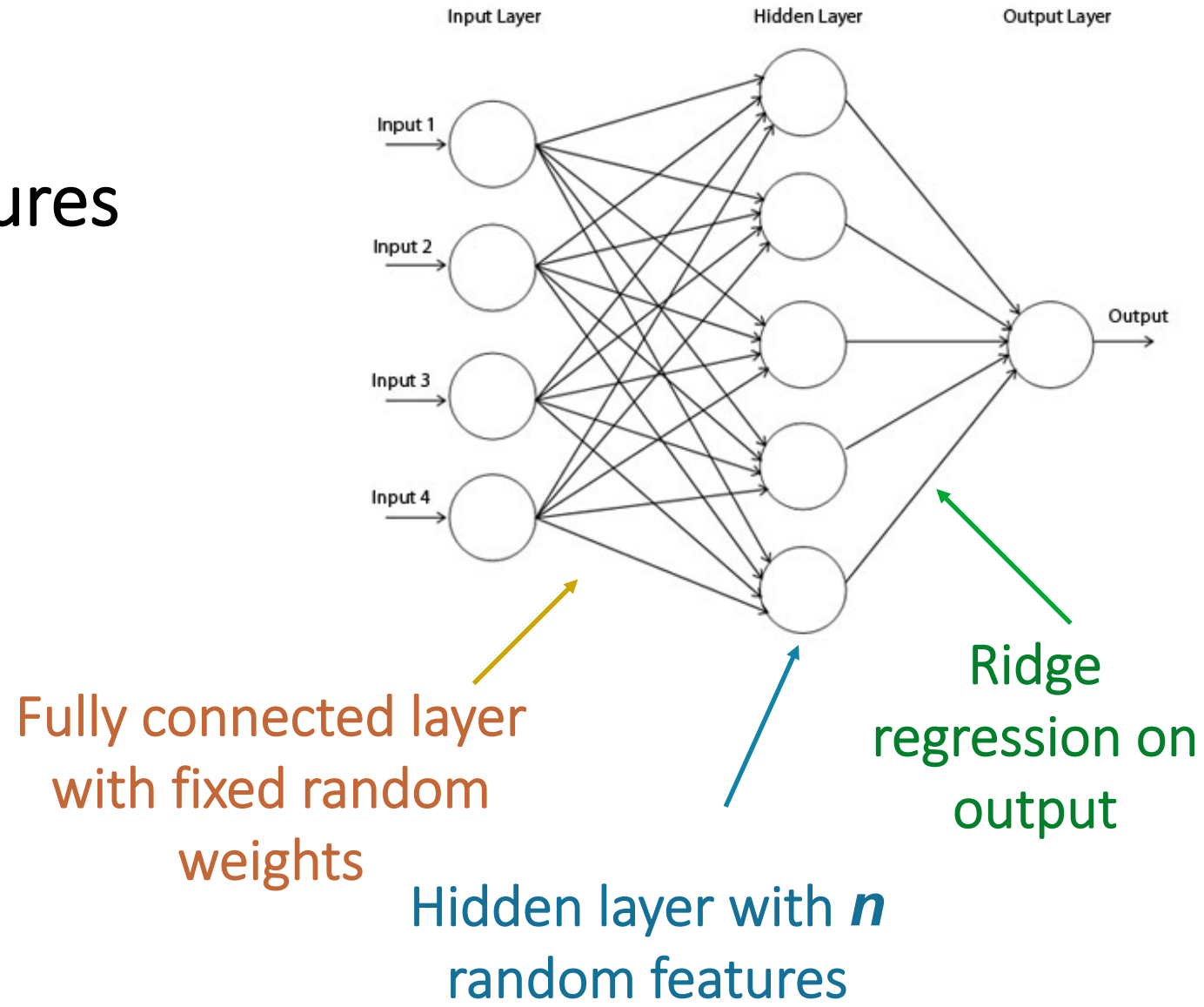




# Case study 1: classification with kernel ridge regression

## Using random features

[in the spirit of Rahimi-Recht]





# Case study 1: classification with kernel ridge regression

---

## Kernel ridge regression

As  $n \rightarrow \infty$ , inner products tend towards a kernel that can be computed explicitly

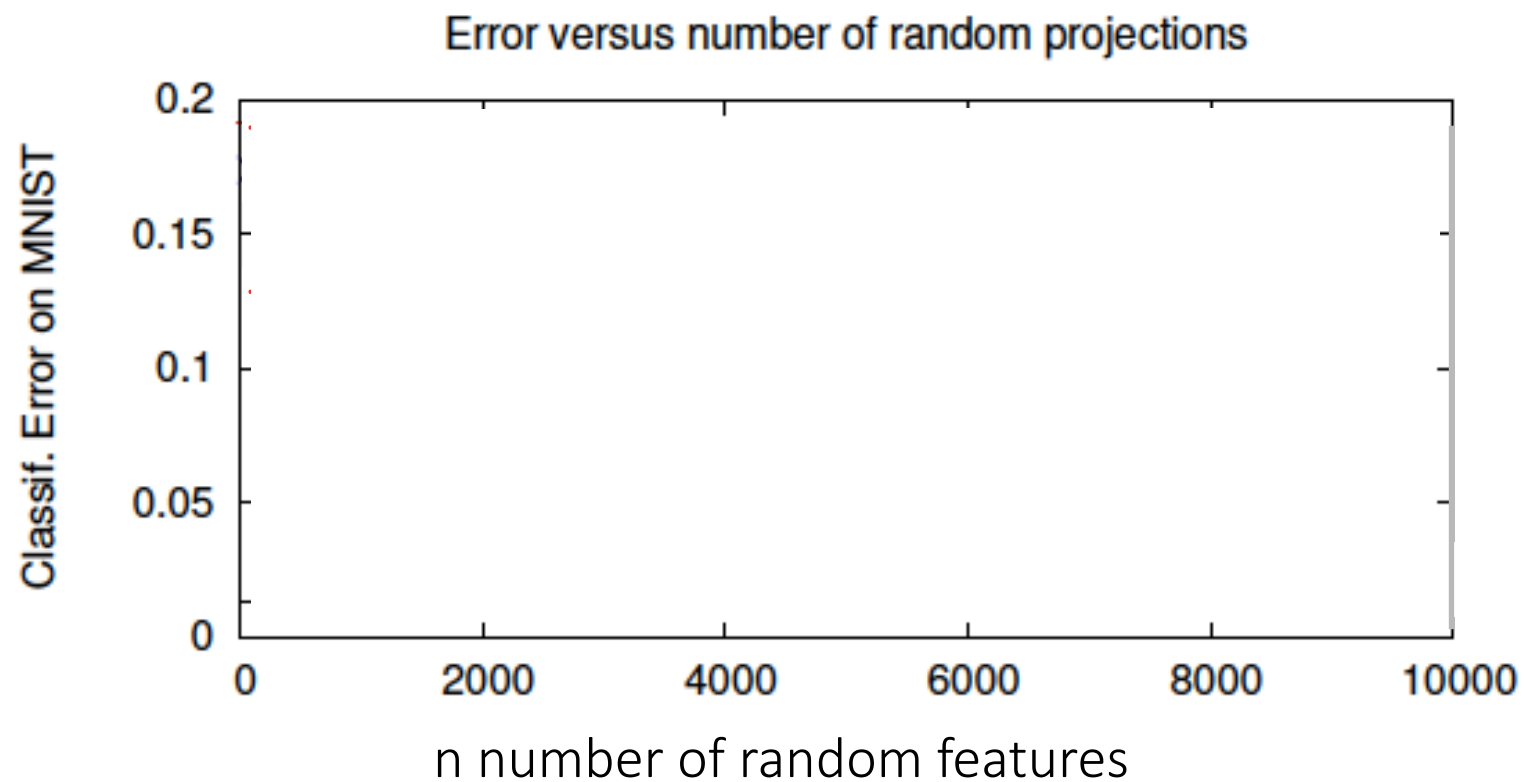
$$k(\mathbf{U}_i, \mathbf{U}_j) = \sqrt{\mathbf{U}_i^T \mathbf{U}_i \mathbf{U}_j^T \mathbf{U}_j} \{2 \mathcal{E}_E[\cos^2 \theta] - \sin^2 \theta \mathcal{E}_K[\cos^2 \theta]\}$$

$\mathcal{E}_K[\cdot]$  and  $\mathcal{E}_E[\cdot]$  are the complete elliptic integrals of the first / second kind  
 $\theta$  is the angle between  $\mathbf{U}_i$  and  $\mathbf{U}_j$

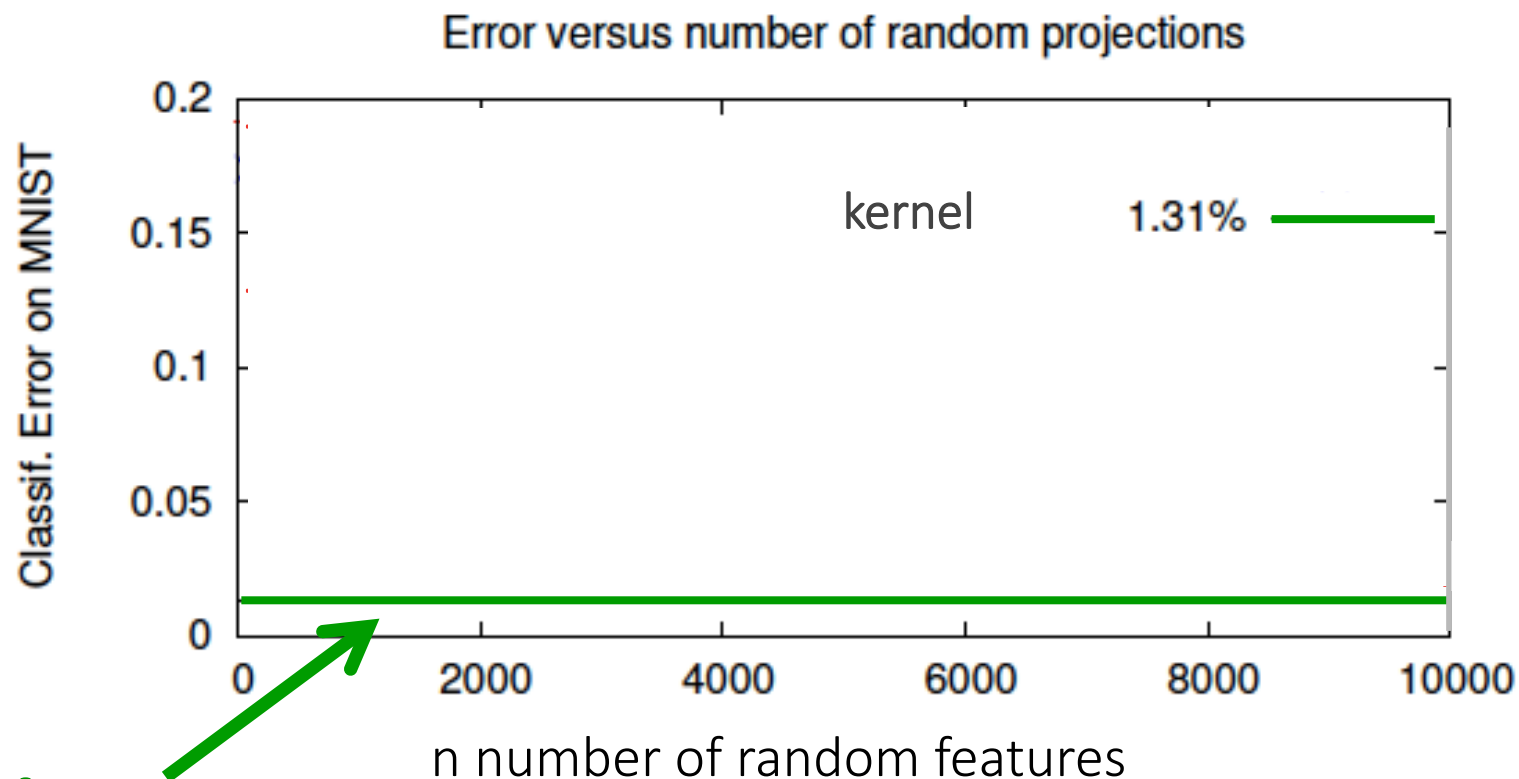
This kernel *numerically* provides a 1.31 % error rate on MNIST

# Case study 1: classification with kernel ridge regression

---



# Case study 1: classification with kernel ridge regression

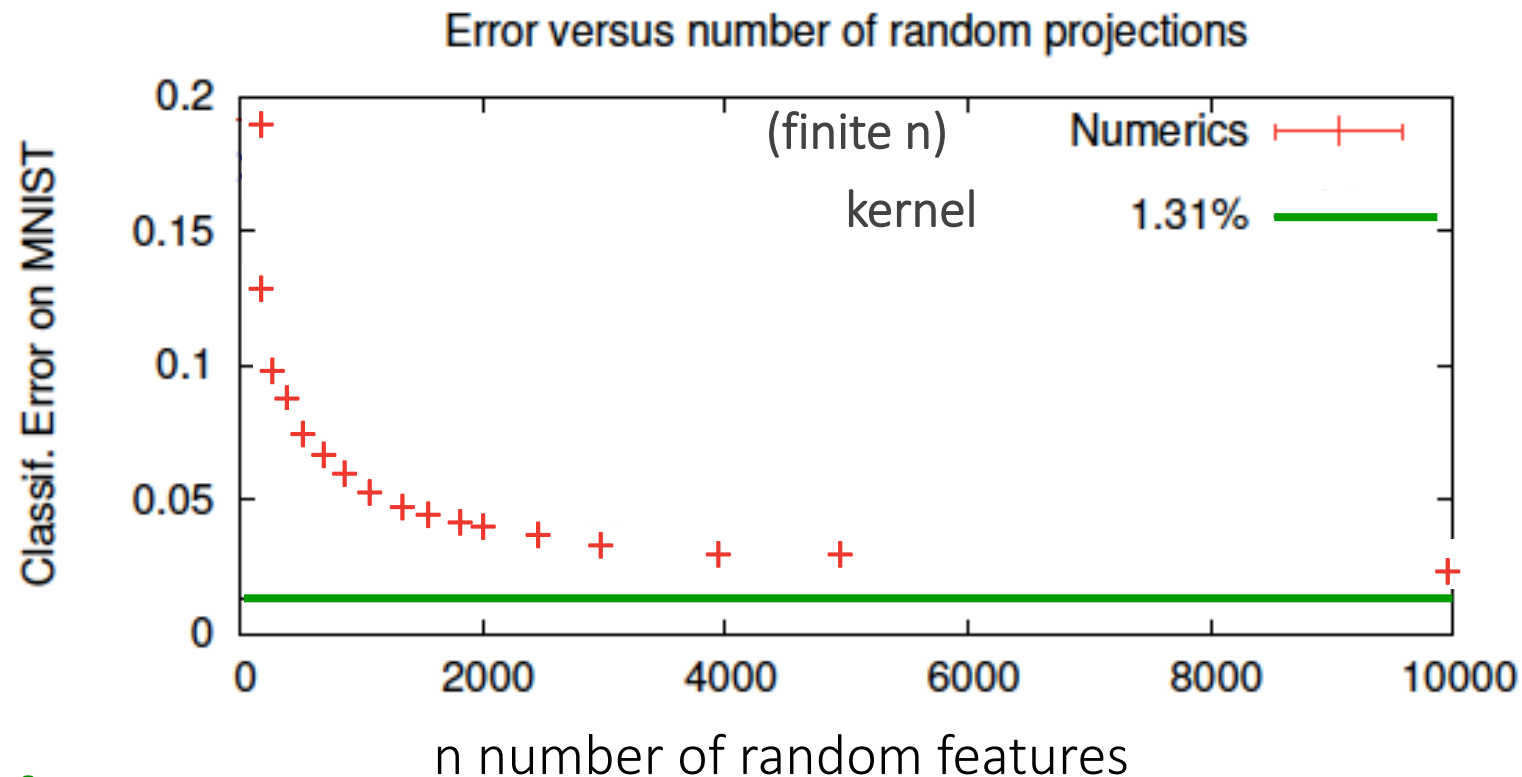


**mathematics**

kernel: asymptotic behavior as  $n \rightarrow \infty$

# Case study 1: classification with kernel ridge regression

## numerical simulations

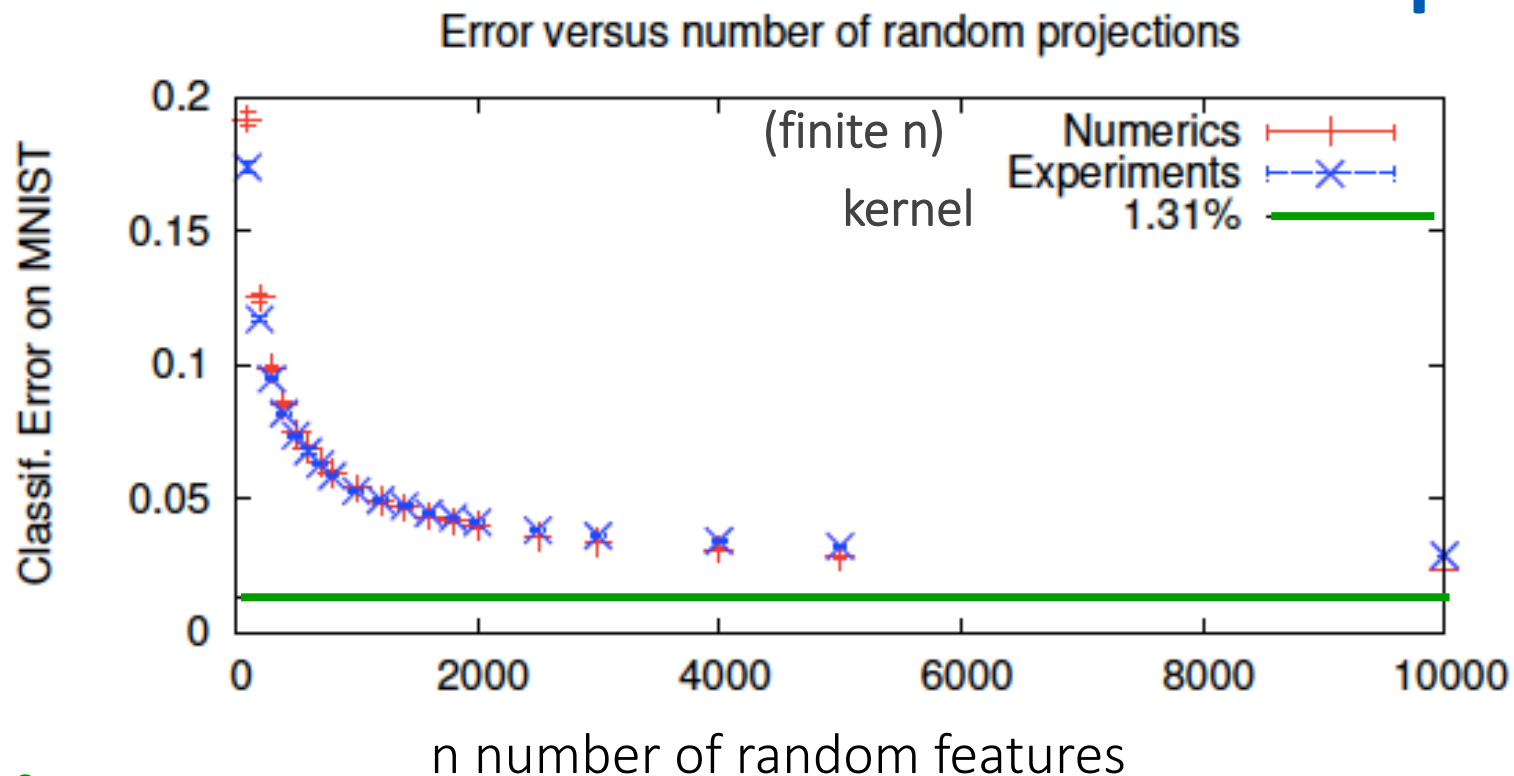


## mathematics

# Case study 1: classification with kernel ridge regression

numerical simulations

optics experiment



mathematics

# Biological motivation for dimensionality expansion



## Science

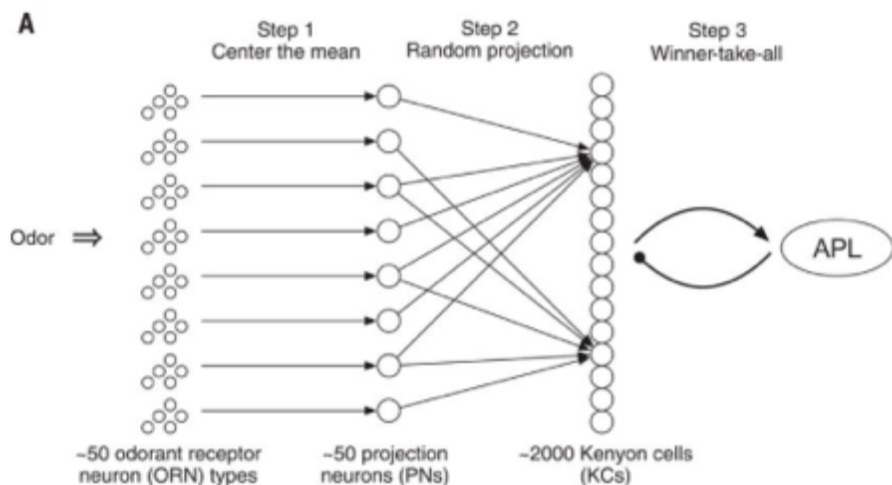
Vol 358, Issue 6364  
10 November 2017

### Fly brain inspires computing algorithm

Flies use an algorithmic neuronal strategy to sense and categorize odors. Dasgupta *et al.* applied insights from the fly system to come up with a solution to a computer science problem. On the basis of the algorithm that flies use to tag an odor and categorize similar ones, the authors generated a new solution to the nearest-neighbor search problem that underlies tasks such as searching for similar images on the web.



[Muhammad M. Karim, GDFL 1.2](#)



## A neural algorithm for a fundamental computing problem

Sanjoy Dasgupta<sup>1</sup>, Charles F. Stevens<sup>2,3</sup>, Saket Navlakha<sup>4,\*</sup>

+ See all authors and affiliations

*Science* 10 Nov 2017:  
Vol. 358, Issue 6364, pp. 793-796  
DOI: 10.1126/science.aam9868



# Case Study 2: Fast Transfer Learning

AI systems « learn » by labelled examples



« cat »

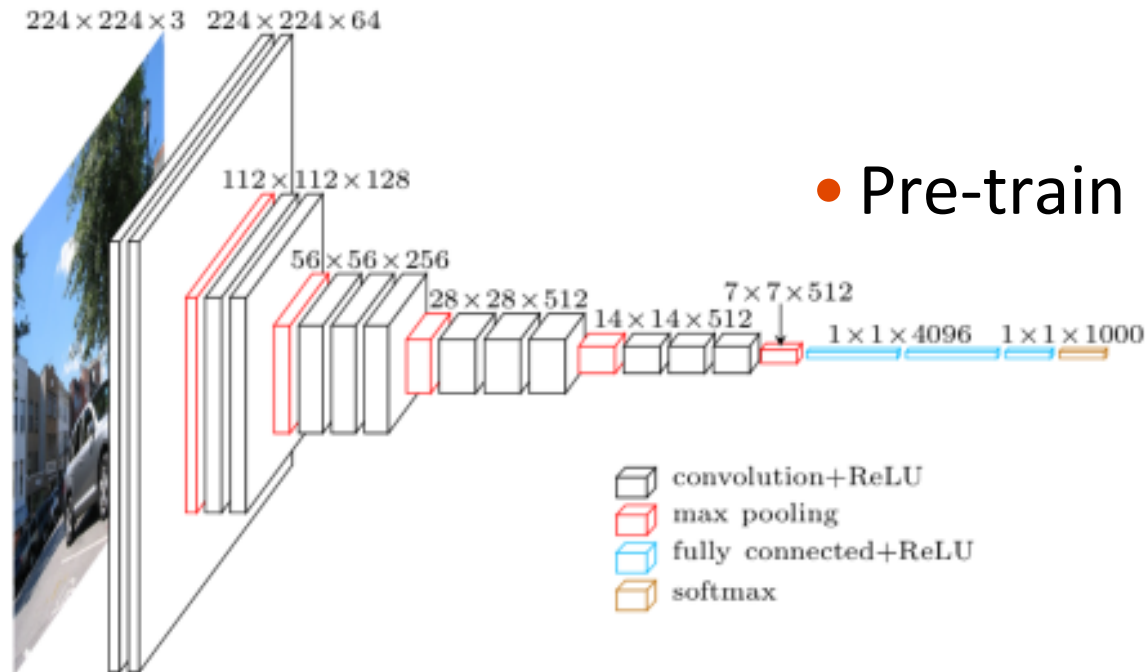


« horse »

- It is an extremely inefficient process: huge amount of data & compute
- « Transfer learning »:
  - pre-train network on large amount of generic data → slow but done once
  - Slightly adapt network to small amount of specific data → fast

# Case Study 2: Fast Transfer Learning

- Start with a standard VGG16 [Simonyan & Zisserman '14] architecture



- Pre-train on generic data with a GPU

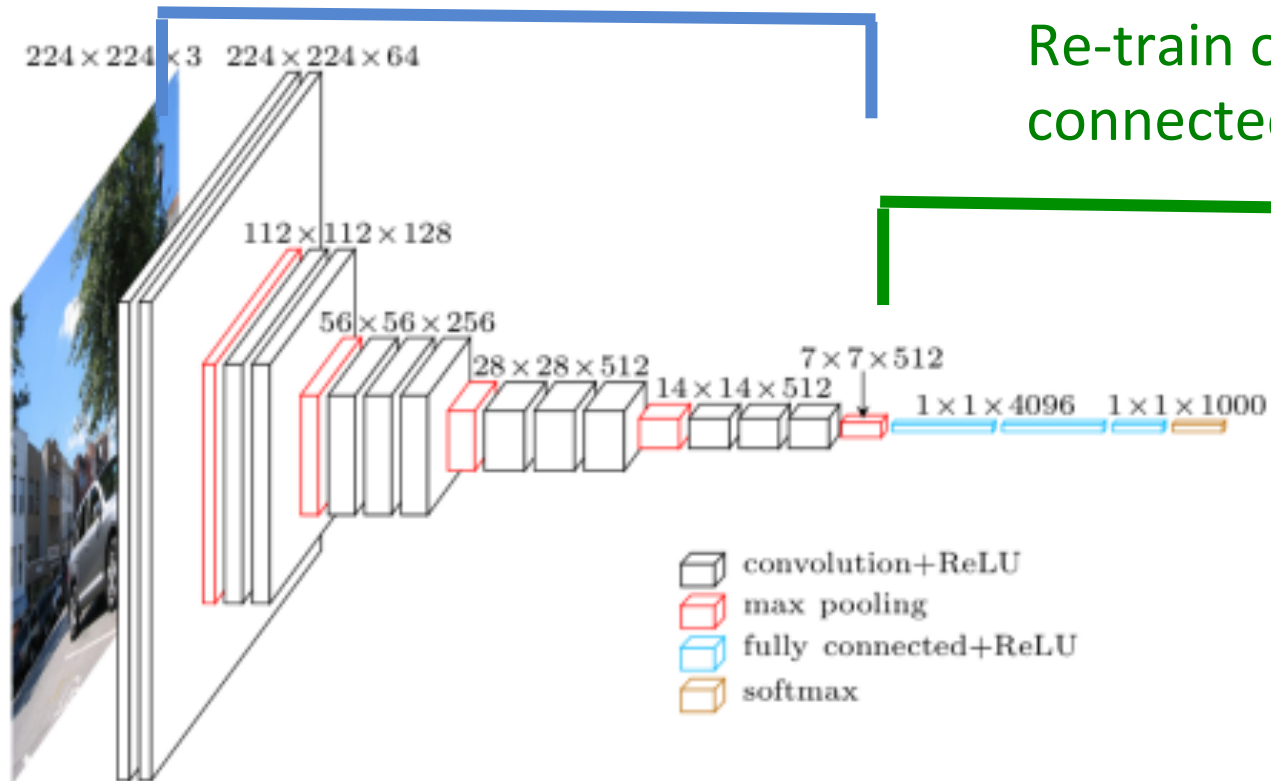


# Case Study 2: Fast Transfer Learning

- Now comes a second dataset : STL10

Keep unchanged

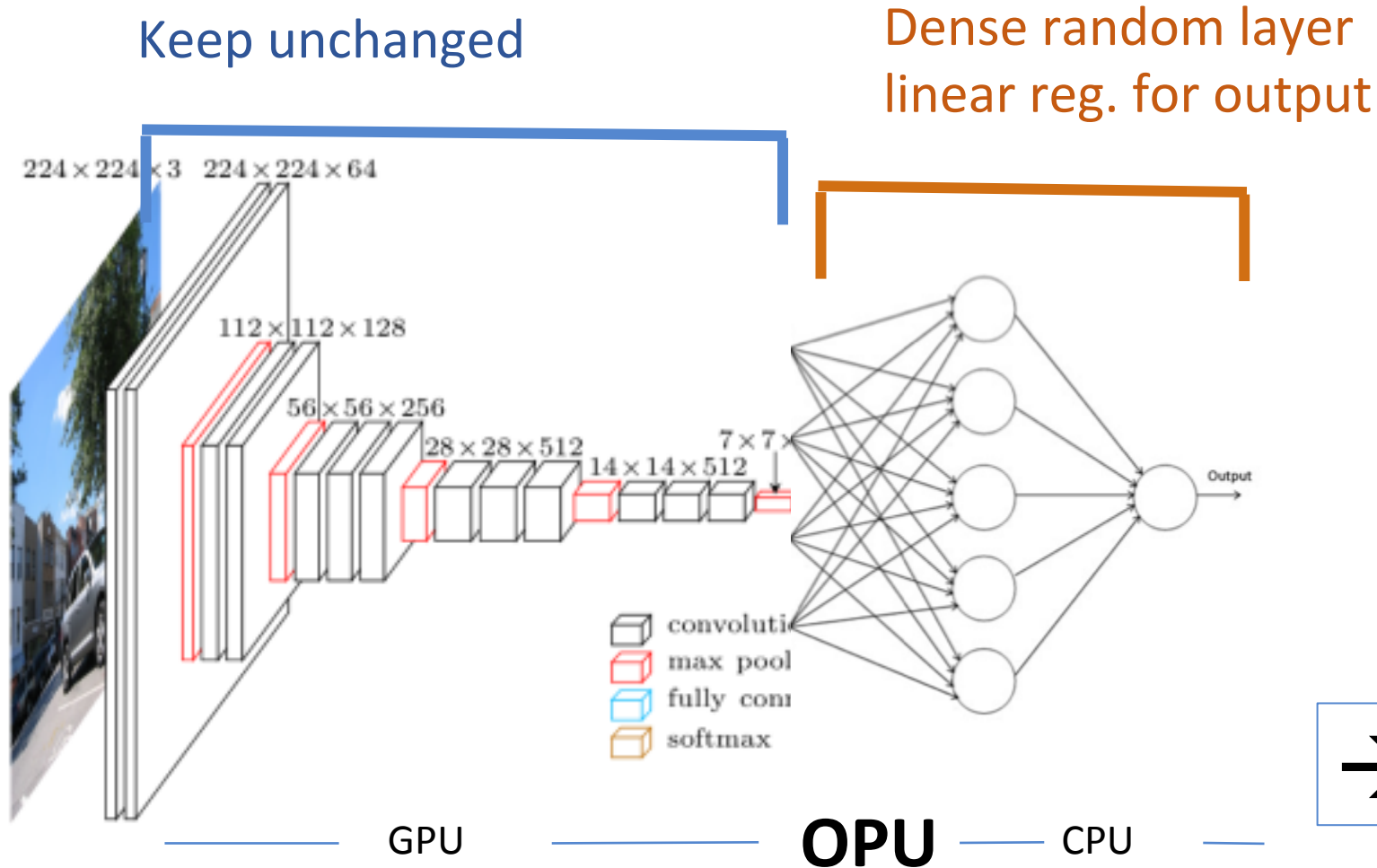
Re-train only fully connected layers



25'

# Case Study 2: Fast Transfer Learning

- Alternative approach



→ Hybrid computing

# Case Study 2: Fast Transfer Learning

[Home](#)

## Transfer Learning Demo

**DATASET** 🐱


Total images	13000
Total classes	10


**OPU RESULTS**


Quantity	OPU	GPU
Test Accuracy [%]	00.0	
Training Time [s]	00.0	
Energy [Wh]	00.0	

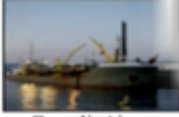
**LightOn**  
We bring Light to AI


vgg16 model loaded successfully.  
Computing convolutional features...


Ground truth ship  
  
Prediction

Ground truth deer  
  
Prediction

Ground truth truck  
  
Prediction

Ground truth ship  
  
Prediction

Ground truth [blurred]  
  
Prediction

Ground truth cat  
  
Prediction

OPU progress

GPU progress

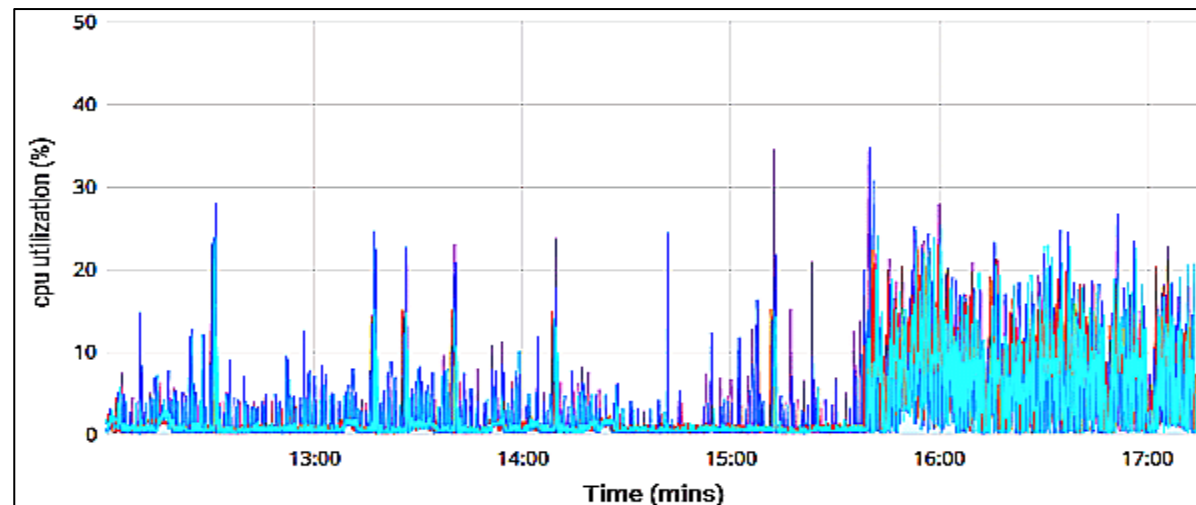
Load model    Convolutions (GPU)    OPU    Classifier (CPU)    Done

NB : x2 video playback speedup

# Case Study 3: Sketching the Change-point

How to detect changes / anomalies on the 1000s of signals monitoring a complex system (factory, datacenter, airplane engine...)?

Example :  
monitoring a  
datacenter



Valsamas *et al*, 2018



# Case Study 3: Sketching the Change-point

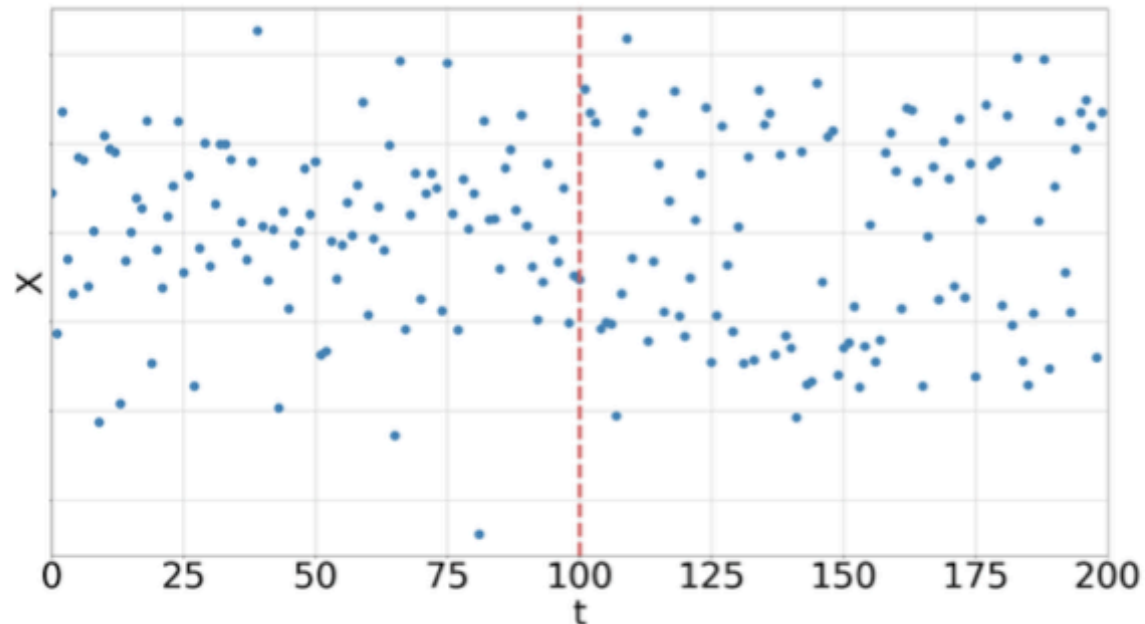
*NEWMA: a new method for scalable model-free online change-point detection*, Keriven, N., Garreau, D., Poli, I., arXiv:1805.08061

Data stream

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t$  with  $x_i \in \chi$

Detect *abrupt* change in

$\mathbb{E}[\Psi(\mathbf{x})]$  with  $\Psi: \chi \rightarrow \mathcal{H}$



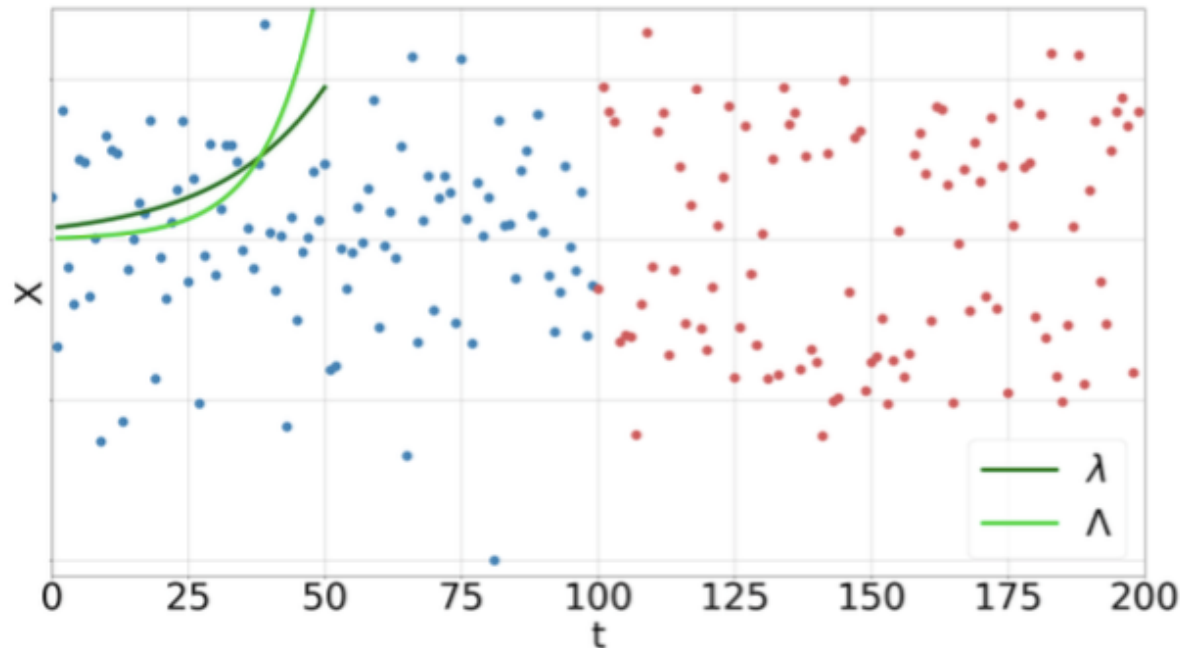
# Case Study 3: Sketching the Change-point

NEWMA equations:

$$\mathbf{z}_1^t = (1 - \lambda) \mathbf{z}_1^{t-1} + \lambda \Psi(\mathbf{x}^t)$$

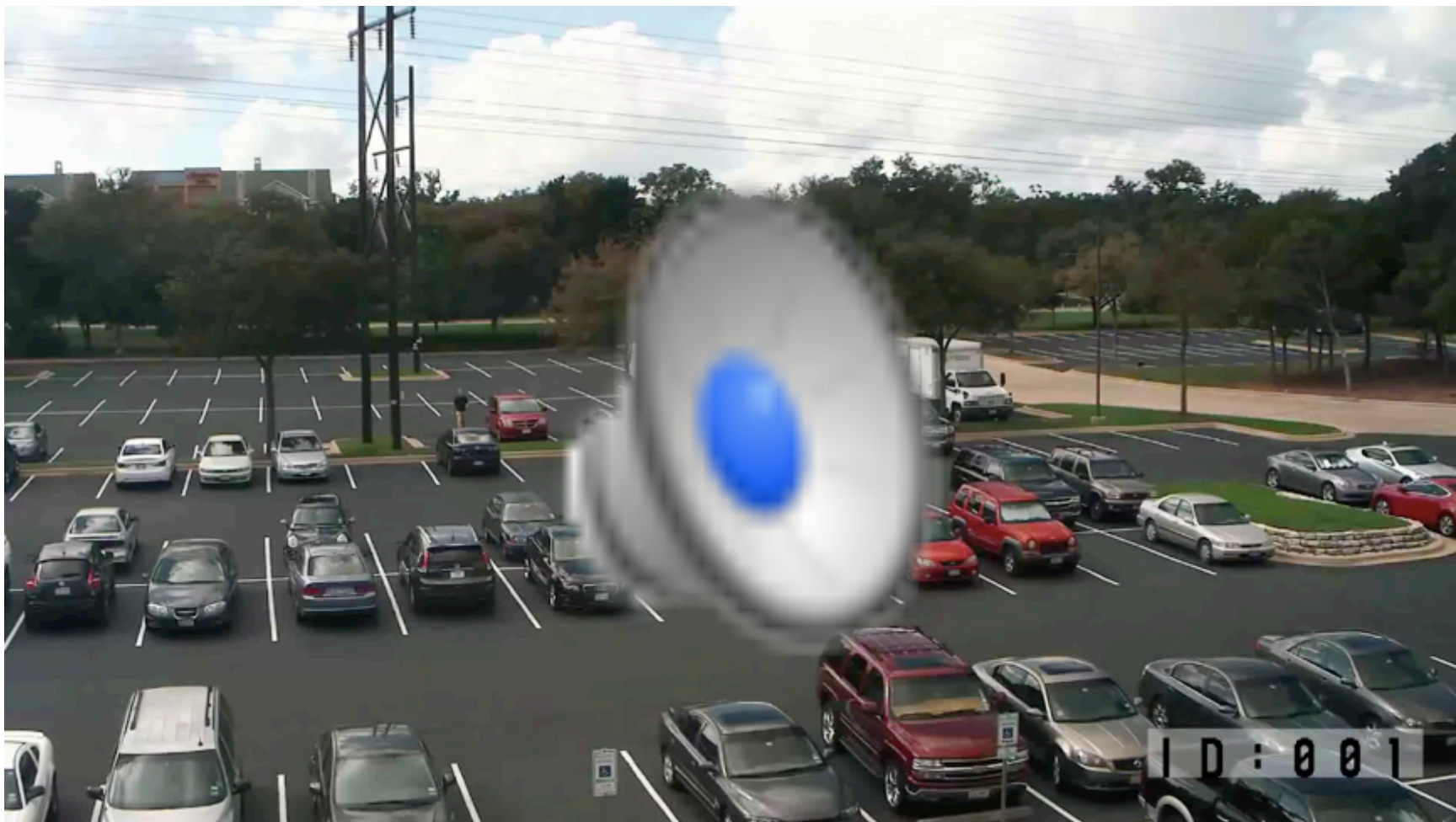
$$\mathbf{z}_2^t = (1 - \Lambda) \mathbf{z}_2^{t-1} + \Lambda \Psi(\mathbf{x}^t)$$

$$0 < \lambda < \Lambda < 1$$

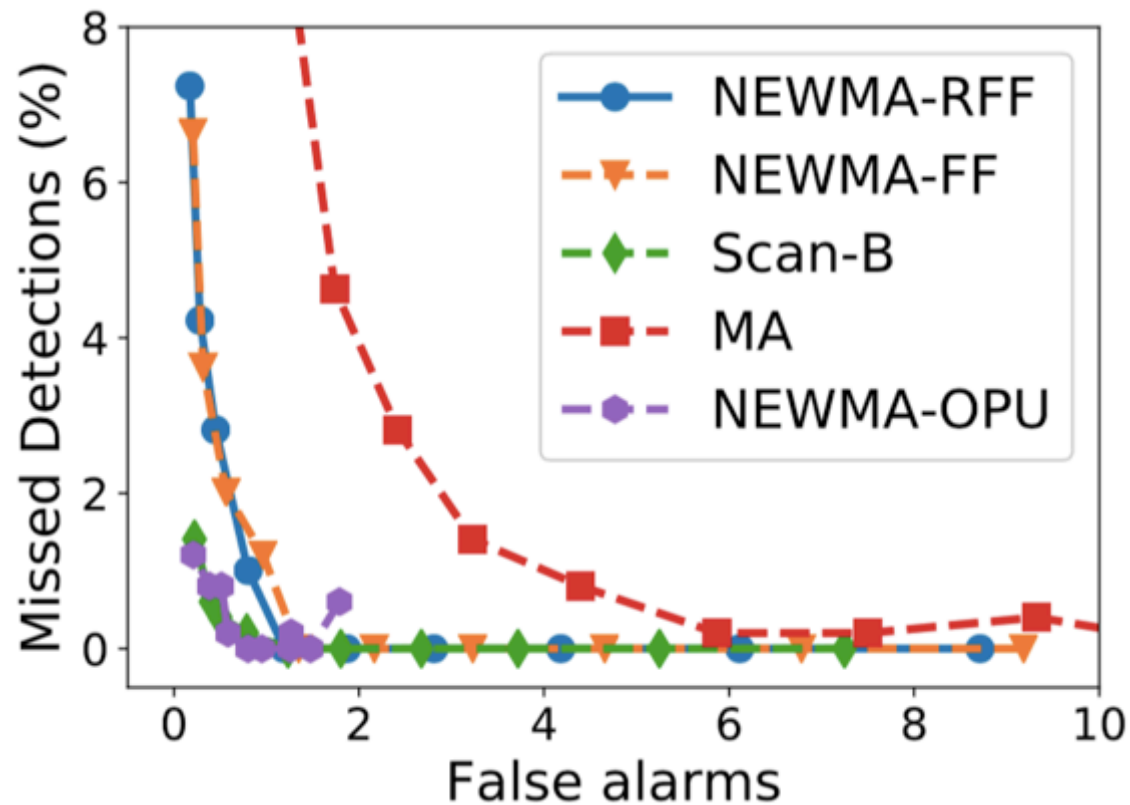
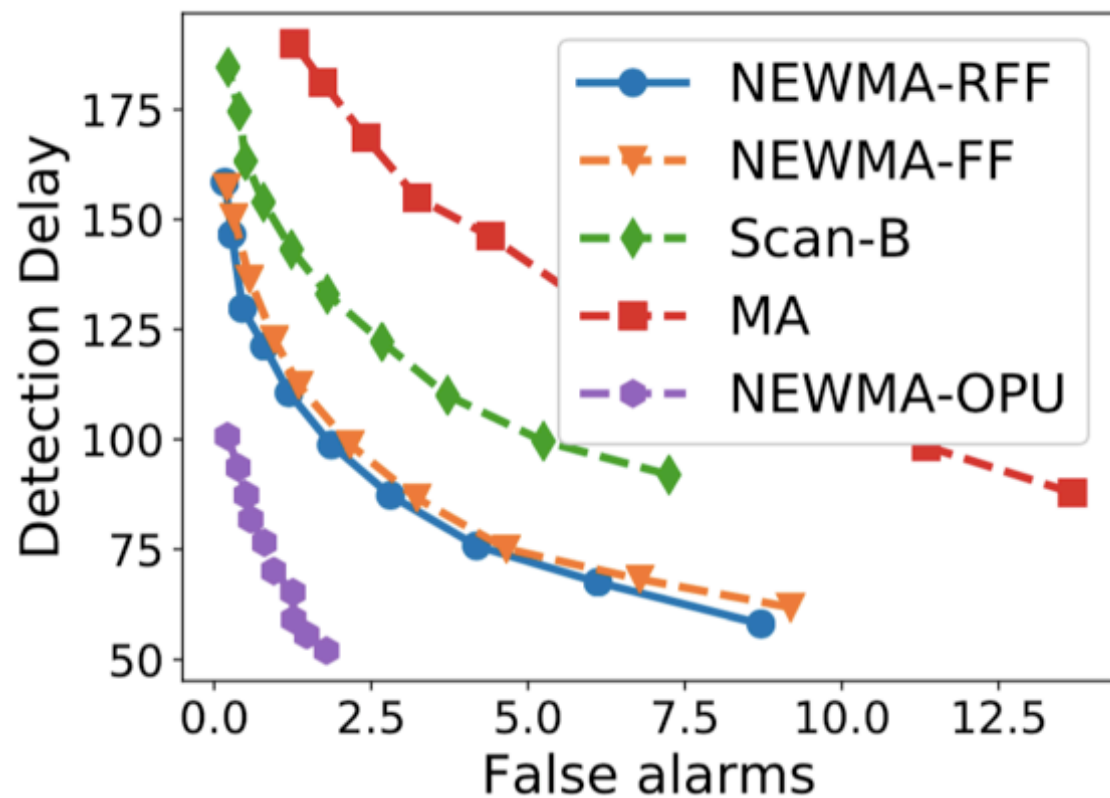


- Compute « sketches » with RPs
- Bounds can be derived
- No need to store samples in memory

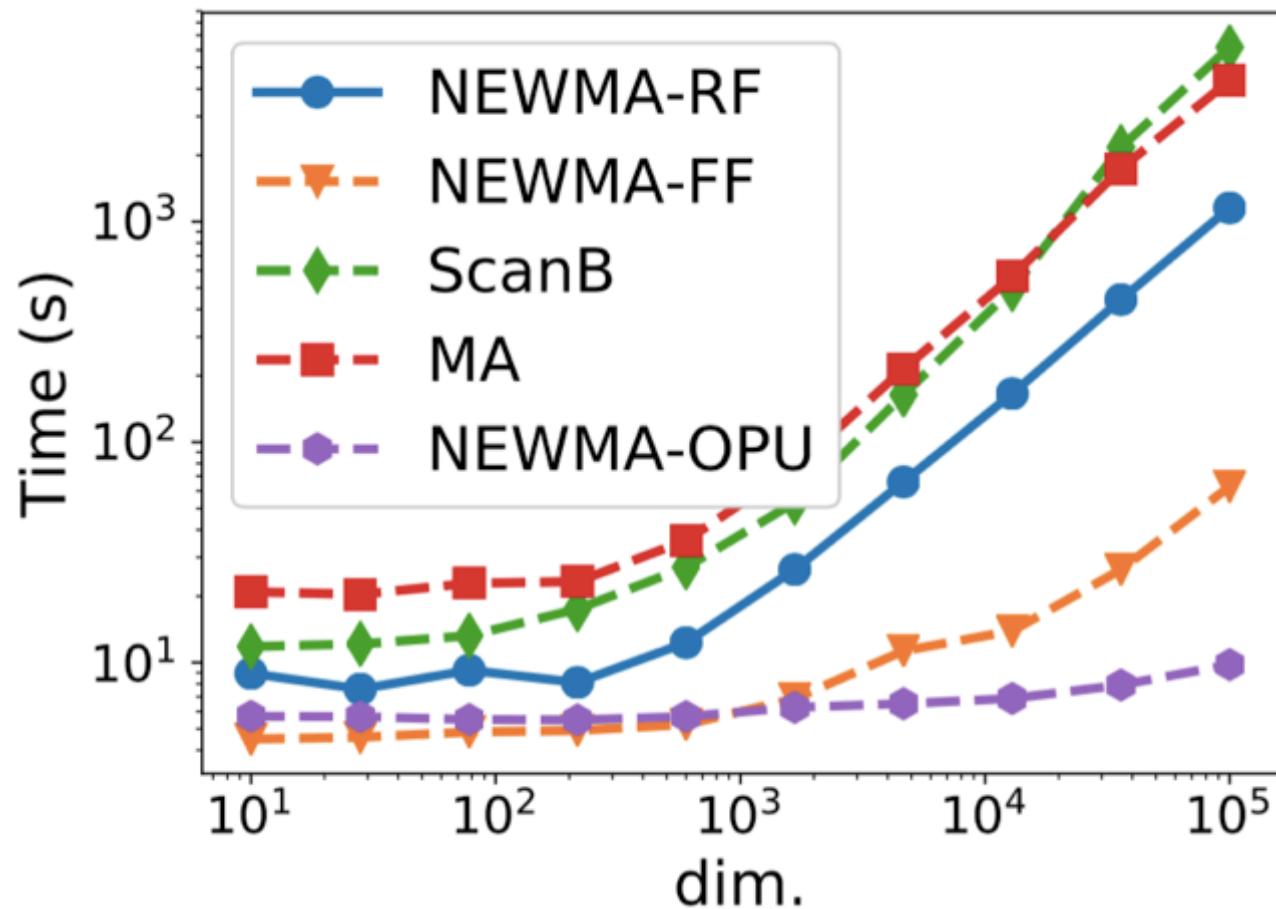
# Case Study 3: Sketching the Change-point



# Case Study 3: Sketching the Change-point

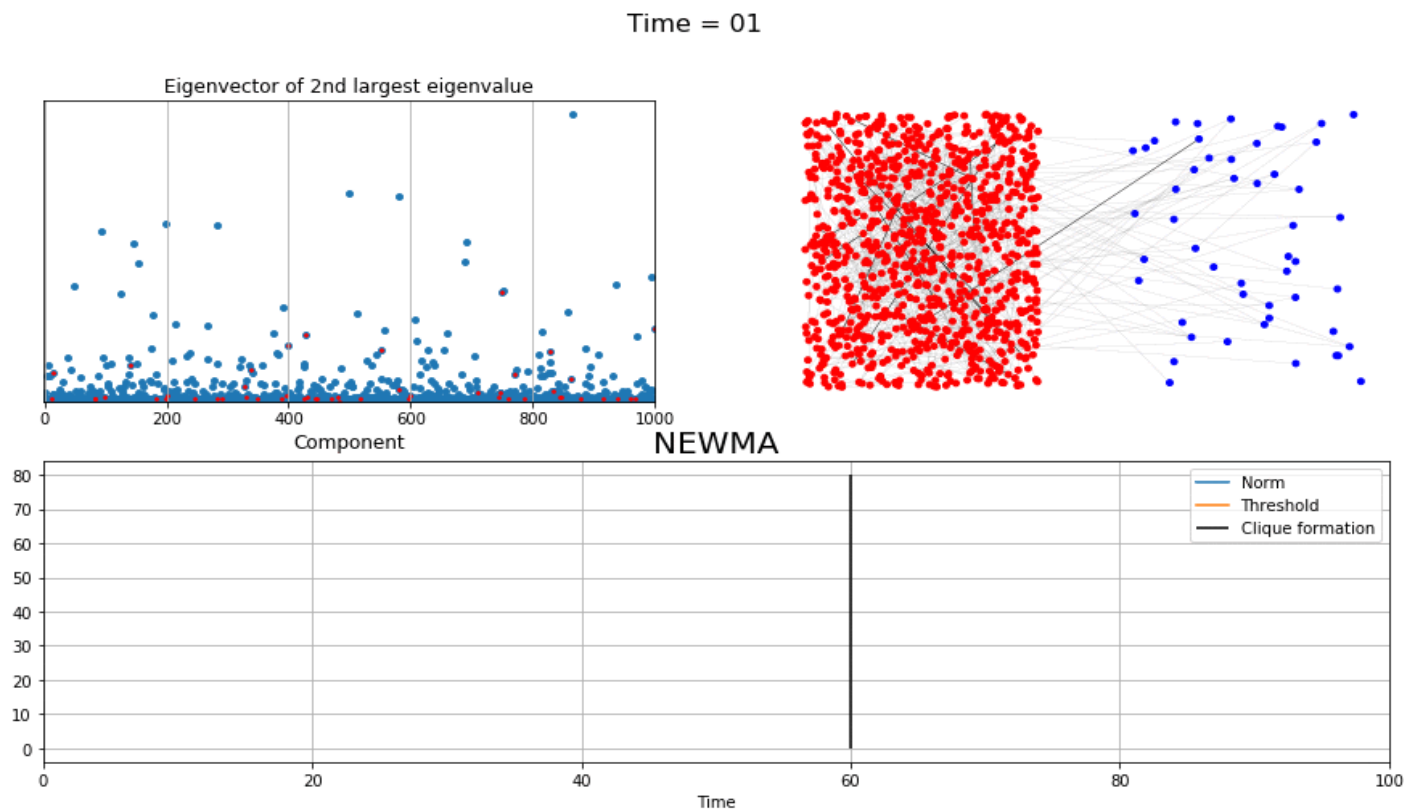


# Case Study 3: Sketching the Change-point



# Case Study 3: Sketching the Change-point

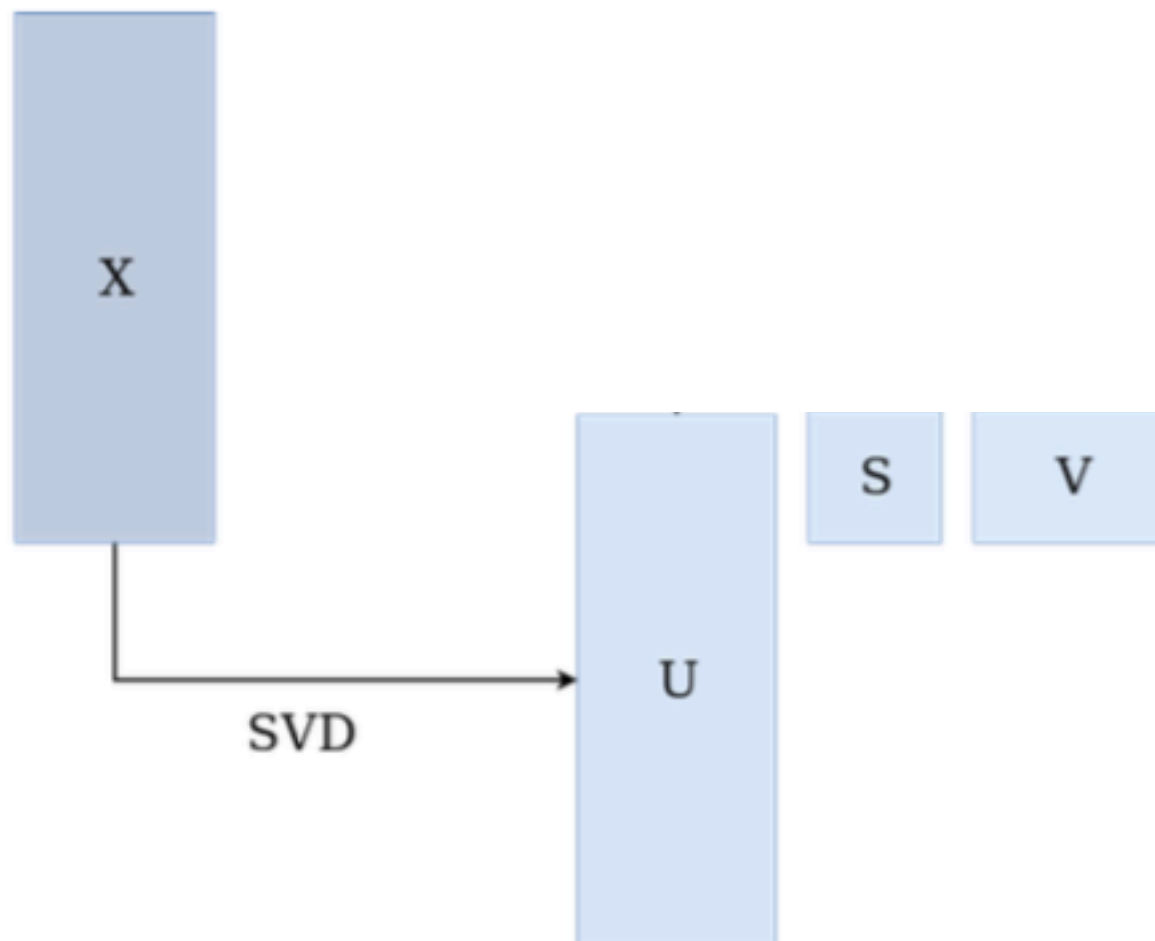
## Clique detection in graphs



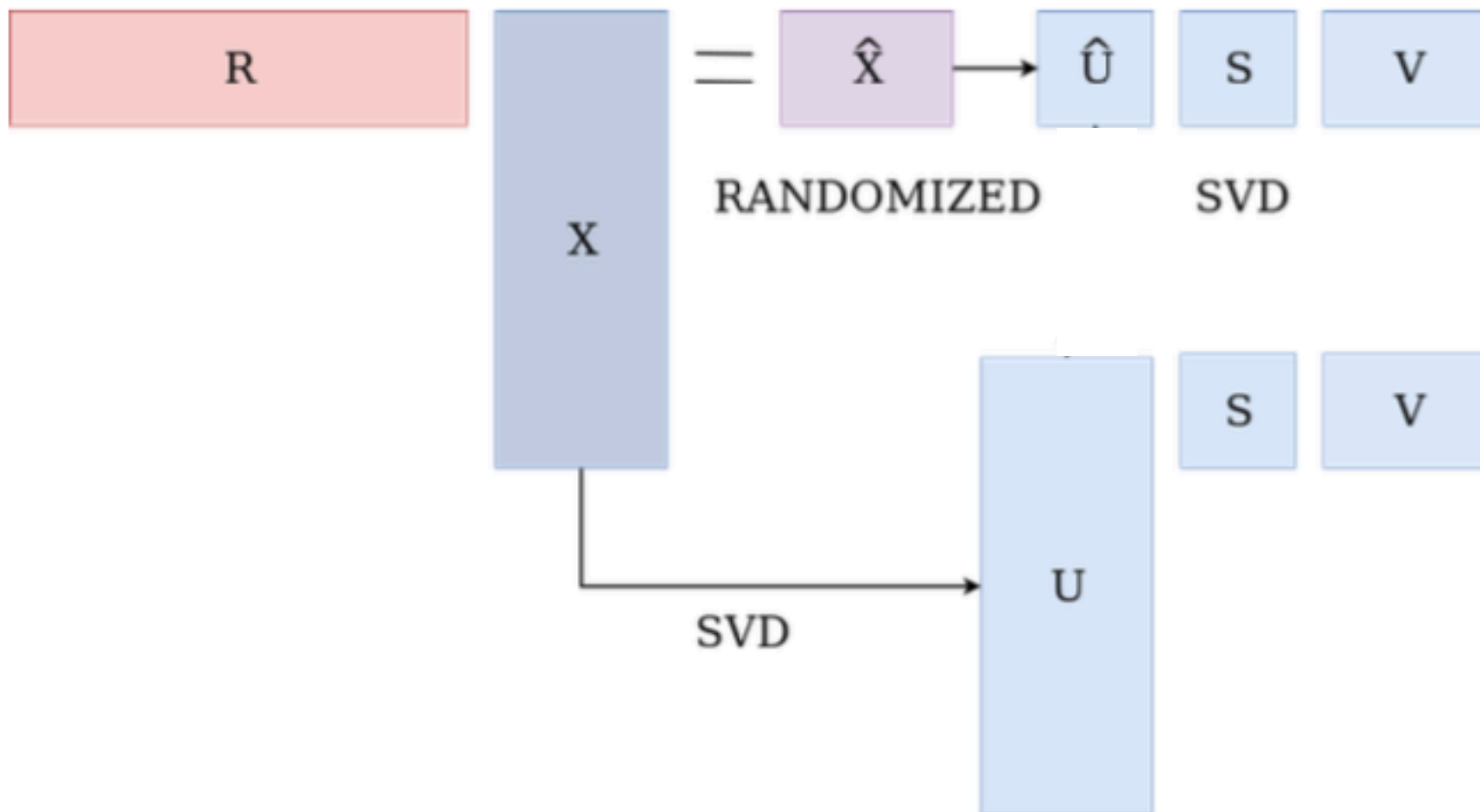


# Case Study 4: Randomized SVD

---



# Case Study 4: Randomized SVD

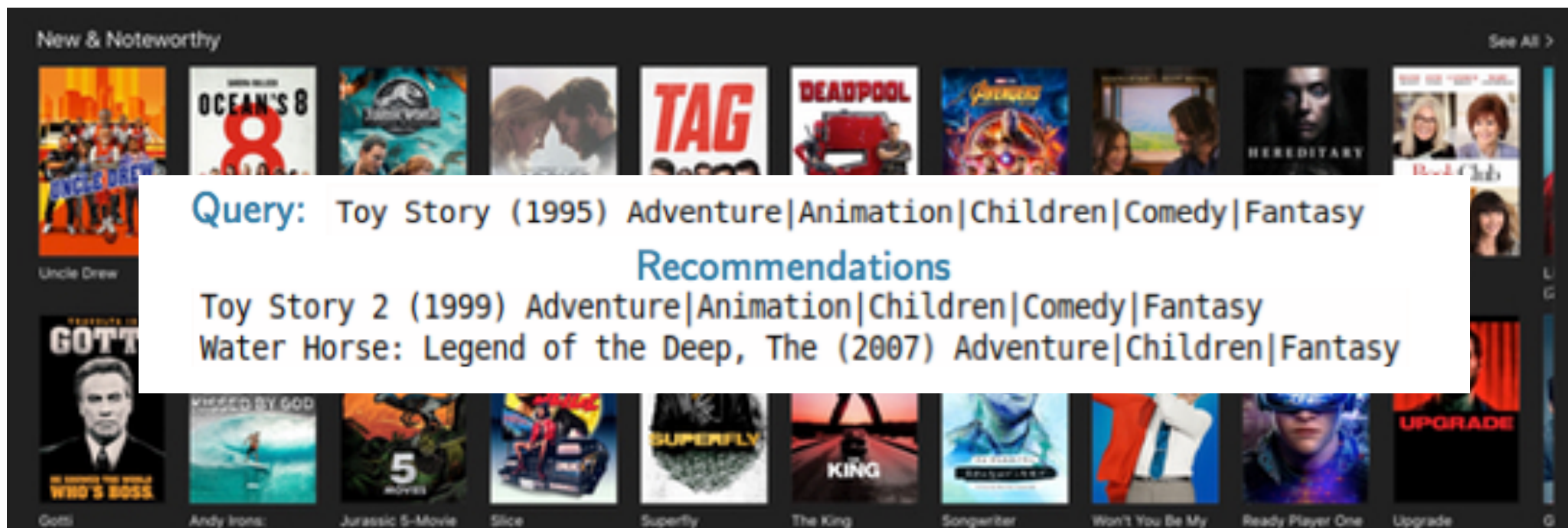


*Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions, Halko, N., Martinsson, P., Tropp, J., 2009, arXiv:0909.4061*

# Case Study 4: Randomized SVD

## Recommender system based on RandSVD

MovieLens 20M database: 27.000 movies x 138.000 users, with 0.5% non-zero entries



The image shows a screenshot of a movie recommendation interface. At the top, there is a header "New & Noteworthy" with a "See All >" link on the right. Below the header is a row of movie posters including "Uncle Drew", "Ocean's 8", "Jurassic World: Fallen Kingdom", "Tag", "Deadpool", "Avengers: Infinity War", "Hereditary", and "Book Club". A white text box is overlaid on the interface, containing the following text:

**Query:** Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy

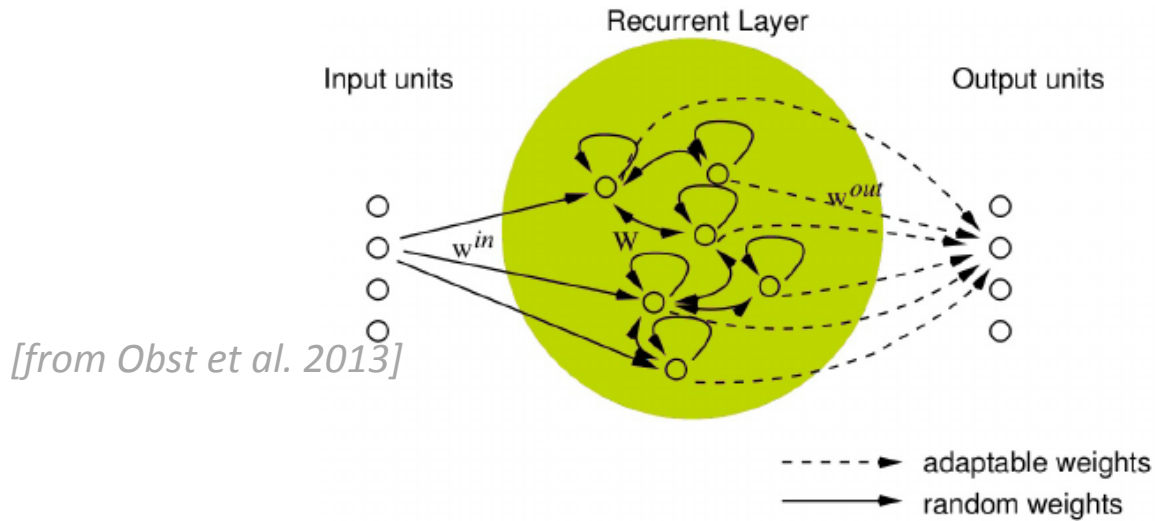
**Recommendations**

- Toy Story 2 (1999) Adventure|Animation|Children|Comedy|Fantasy
- Water Horse: Legend of the Deep, The (2007) Adventure|Children|Fantasy

Below the text box is another row of movie posters including "Gotti", "Kissed by God", "Jurassic 5-Movie", "Slice", "Superfly", "The King", "Songwriter", "Won't You Be My", "Ready Player One", and "Upgrade".

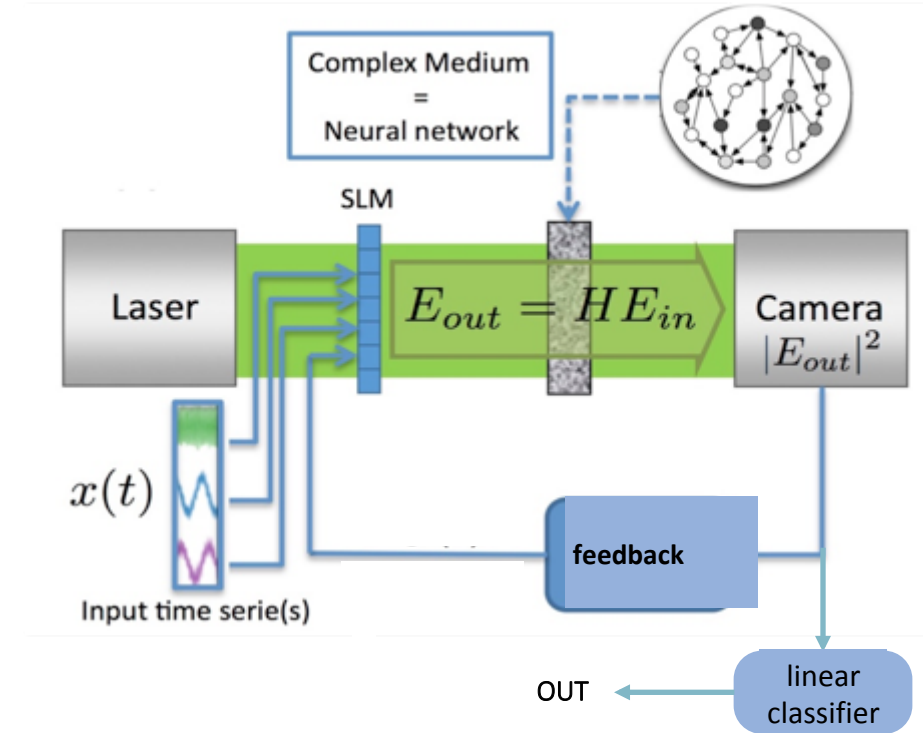
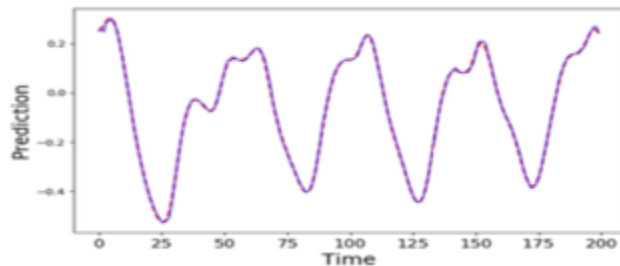
# Case Study 5: Optical Echo-State Network

A physical implementation of large-scale echo-state networks



regression on complex time series

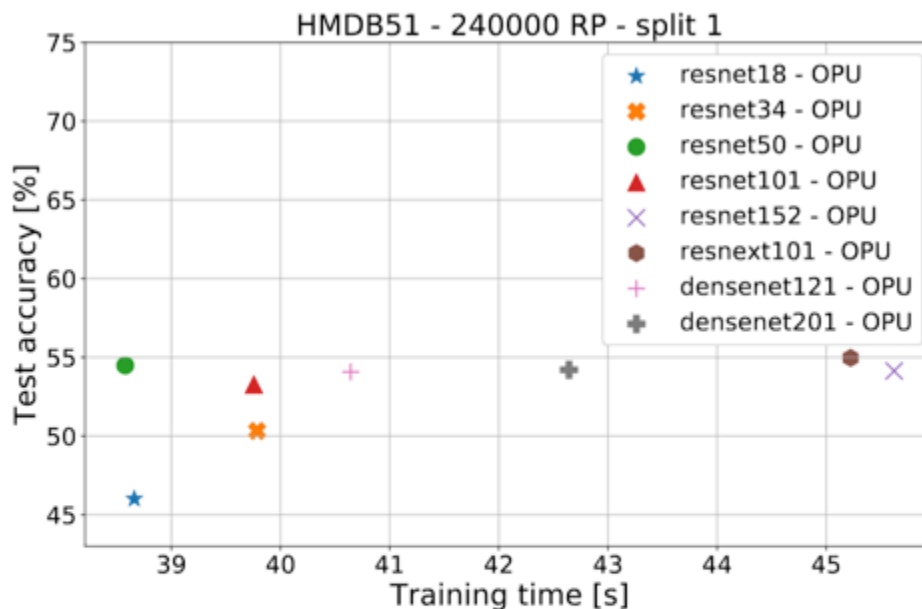
Ex: predict dynamics of Mackey-Glass eqs. (Dong. et al)



2 orders of magnitude larger than standard PCs

# Case Study 6: Video classification

Transfer Learning on videos [1] using the Kinetics dataset



Method	HMDB-51
ResNet-18 (scratch)	17.1
ResNet-18	56.4
ResNet-34	59.1
ResNet-50	61.0
ResNet-101	61.7
ResNet-152	62.4
ResNet-200	63.5
DenseNet-121	59.6
ResNeXt-101	<b>63.8</b>

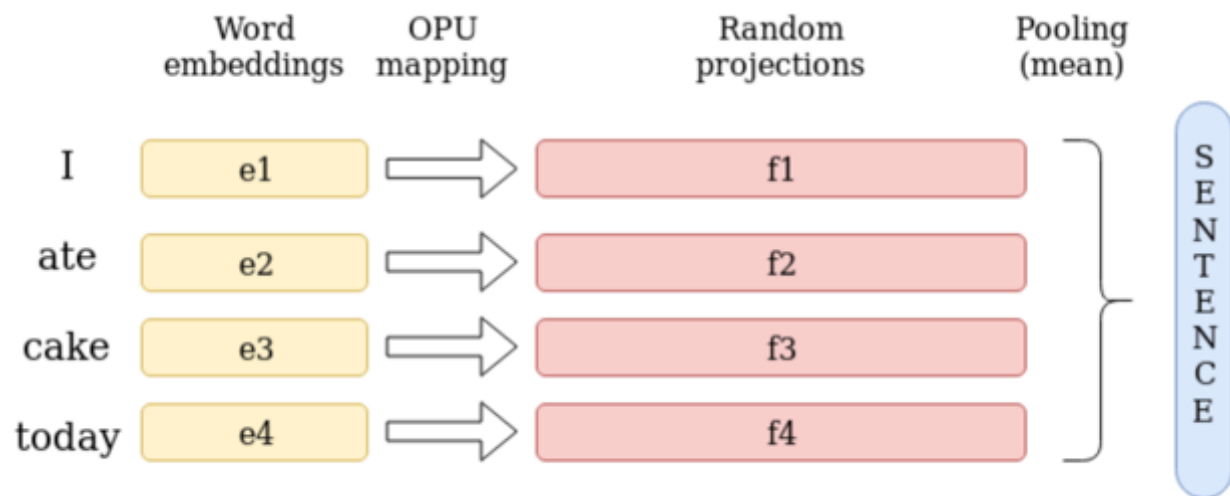
**Left:** Accuracy versus training time for the HMDB51 dataset on different pretrained Kinetics models. **Right:** Accuracy with backprop.

Standard training times for 3D-CNN:

- 150/200 epochs with 8 GPUs [2]
- 5k epochs "using just 16 GPUs" [3]

[1] Hara, Kensho, Hirokatsu Kataoka, and Yutaka Satoh. "Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?." 2018  
 [2] DIBA, Ali, et al. Temporal 3d convnets: New architecture and transfer learning for video classification. arXiv preprint arXiv:1711.08200, 2017.  
 [3] Carreira, Joao, and Andrew Zisserman. "Quo vadis, action recognition? a new model and the kinetics dataset." 2017.

# Case Study 7: Natural Language Processing



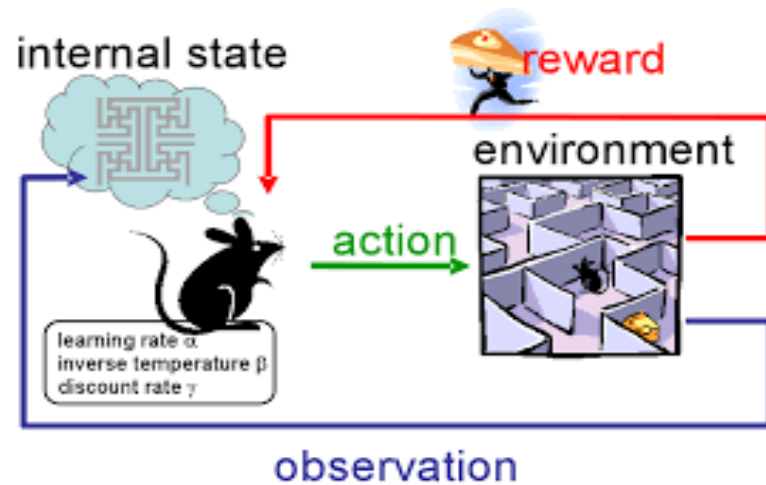
	Ours	SkipThought	InferSent
MR	79.4	79.4	<b>81.1</b>
CR	81.6	83.1	<b>86.3</b>
MPQA	87.7	89.3	<b>90.2</b>
SST2	82.7	82.9	<b>84.6</b>
SUBJ	92.9	<b>93.7</b>	92.4
TREC	<b>89.2</b>	88.4	88.2
SICK-E	83.5	79.5	<b>86.3</b>
SICK-R	84.3	85.8	<b>88.3</b>
MRPC	75.2	73.2	<b>76.2</b>
STSB	70.6	68.9	<b>75.6</b>

Table 1: Results on downstream tasks of the SentEval toolkit [2]. Our model (15,000 random projections) is inspired from [1] and gets comparable results to much more sophisticated works. SkipThought [3] took a month to train and InferSent [2] requires annotated data that is thus far available only in English.



# Case study 8: Reinforcement learning

How to learn complex tasks through feedback ?



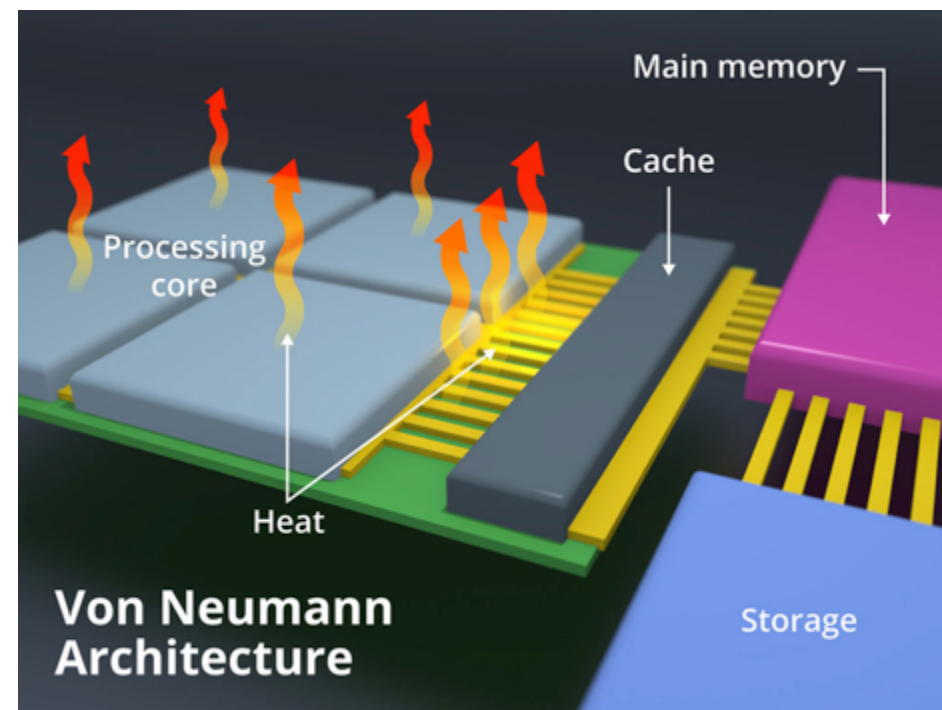
Becominghuman.ai

# Case study 8: Reinforcement learning

This system uses the OPU operation to search in past sequences a similar position

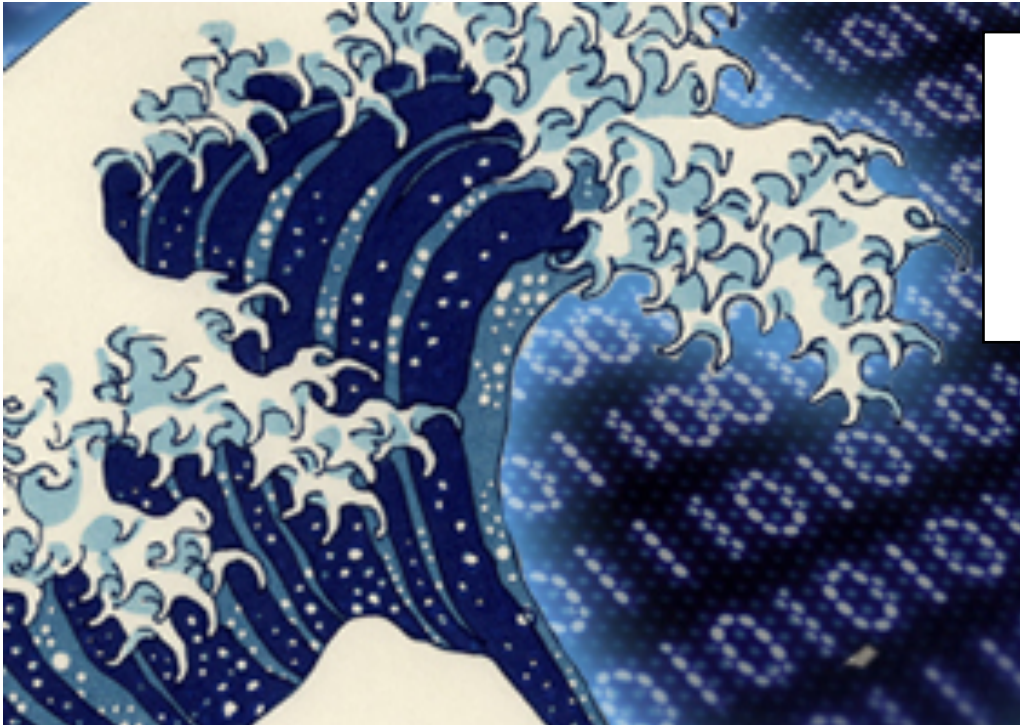


Will Von Neumann architectures  
stay prevalent in the AI era ?



<http://www.rochester.edu/newscenter/microprocessors-computing-architecture-304252/vonneumann-architecture/>

# A conclusion and an invite



Optical Co-processing already allows  
some ML computations **at scale**  
**Fast ML prototyping**

- Kernel Ridge Regression
- Fast Transfer Learning
- Echo-State Networks
- Change-point detection in streams and graphs
- Recommender systems with RandSVD
- ...

**What's your case study ?**  
laurent@LightOn.io

<https://blog.apterainc.com/the-data-tsunami-is-coming.-is-your-company-ready>

# Selected references

---

- "Imaging With Nature: Compressive Imaging Using a Multiply Scattering Medium", A. Liutkus et al., Scientific Reports 4 (july 2014)
- "Reference-less measurement of the transmission matrix of a highly scattering material using a DMD and phase retrieval techniques", A. Drémeau et al., Optics Express 23(9), 2015
- "Random Projections through multiple optical scattering: Approximating kernels at the speed of light", A. Saade et al., Proc. ICASSP (2016)
- "Scaling up Echo-State Networks with multiple light scattering", J. Dong et al., arXiv:1609.05204
- "NEWMA: a new method for scalable model-free online change-point detection", Nicolas Keriven, Damien Garreau, Iacopo Poli, arXiv:1805.08061
- "Machine Learning and the Physical Sciences », G. Carleo et al., arXiv 1903.10563
- "Principled training of Neural Networks with Direct Feedback Alignment", J. Launay et al., arXiv 1906.04554