# Deep Learning for AI: turning the GPS on

Cesare Furlanello

FBK // DATASCIENCE // MPBA

**@furlanello**

mpbalab.fbk.eu

Smart Cities & Communities, Health & Well Being, Future Media, Machine Translation, Smart Digital Industry, Data Science, CyberSecurity

FBK for Artificial Intelligence – 2018

# REVOLUTIONARY CHANGE OF LANDSCAPE IN RESEARCH & INNOVATION

WORLD COMPETITION

RISKS FOR SOCIETY

HYPES

FREE SPACE

RESOURCES

**Private traits and attributes are predictable from digital records of human behavior**

**IOT, Industry, Retail, Finance, <u>Healthcare</u>, Agricolture already changed**

**Machine Learning as commodity (0.9 $/h)**

**0.3-0.5% of the population knows how to code**

# WHERE TO START ?

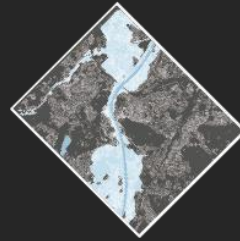## PRODUCTS

OBSERVATION

FORECAST
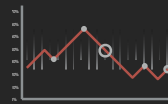
INTEGRATION OF EXISTING SYSTEMS

UNMANNED AI SOLUTIONS

DECISION SUPPORT SYSTEMS

## BIG / DATA

DATA LANDSCAPE

ORGANIZATION PROCESS

HISTORICAL DATA

SOCIAL MEDIA DATA STREAMS

## HUMAN RESOURCES

AI / ML PRACTITIONERS

?

## TECHNOLOGIES

SMARTPHONE APP

WEB INTERFACES

WEARABLES

*C. Furlanello – MPBA Nov2017*

Credits:
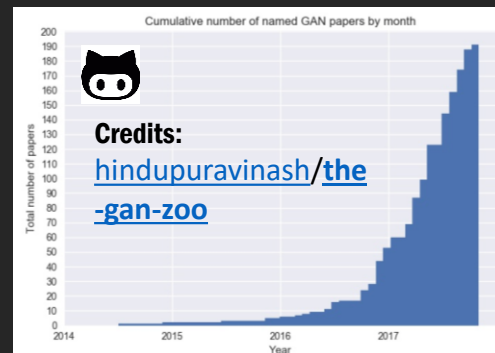**I-REACT H2020**

# UNDER THE HOOD



**"It's ML, not magic"**  Credits: @smerity

- **The Data Science stack (Python, R) and resources ( arXiv.org )**
- **Keras with TensorFlow backend**
- **PyTorch**
- **Fast, well tuned baselines**
- **Model Selection: human intuition in Deep Learning is bad**
- **Ability to accurately measure progress over time**
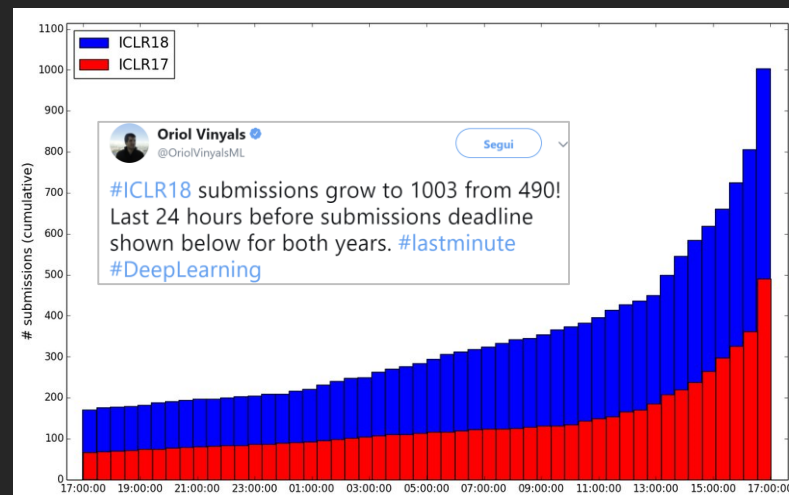
**OK, but … do I need a GPU Armada?**

# FLOOD & FIRE

Cumulative number of named GAN papers by month

**Credits:**
hindupuravinash/**the -gan-zoo**

## unprecedented and growing cost of keeping track of new DL ideas

- Karpathy's http://www.arxiv-sanity.com/
- http://search.iclr2018.smerity.com/:
  hits RNN: 324, CNN: 455, GAN: 666



Oriol Vinyals @OriolVinyalsML

#ICLR18 submissions grow to 1003 from 490! Last 24 hours before submissions deadline shown below for both years. #lastminute #DeepLearning

# MPBA LAB

**Data Science Lab (25 pe) : Maths, CS, BioEngineering, Physics, …**

Unifarm's robotic warehouse

- **Machine Learning: systems that can learn from examples and predict over novel data**

- **Applied: Predictive Biomarkers, agri-tech; environmental risk, IoT, personal sensors**

**2 ½ STARTUPS**

- **MPA Solutions: geospatial analytics**

- **Motorialab: sport analytics ➔ insurtech**

- **Multipl.AI: DL for Precision Medicine**

**ML Projects with industry and retail**

**EIT Digital: Wearable Analytics, Smart Retail, Fraud profiling**

# MPBA LAB

- **ML as infrastructure: Big Data Analytics, Networks, GIS, bio-informatics, <u>Deep Learning</u>**

- **Fast deployement of <u>dashboard analytics and microservices</u> (CPU/GPU) in cloud**

**EXCELLENCE IN RESEARCH**

- **<u>Precision Medicine</u>, with FDA, OPBG, Riken; HEP with Liverpool**

**CHALLENGES FOR DEEP LEARNING**

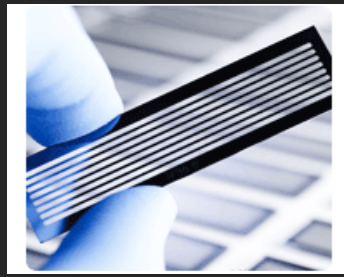- **Data integration & reproducibility in diagnostics and pharmacogenomics**



Forecasting EDM Ticket Sales with ML & Networks

**Dashboards as enablers of "actionable" analytics: explore patterns, interact with models for simulation**
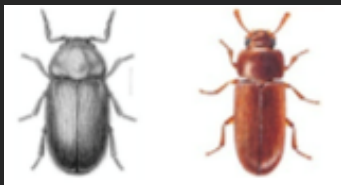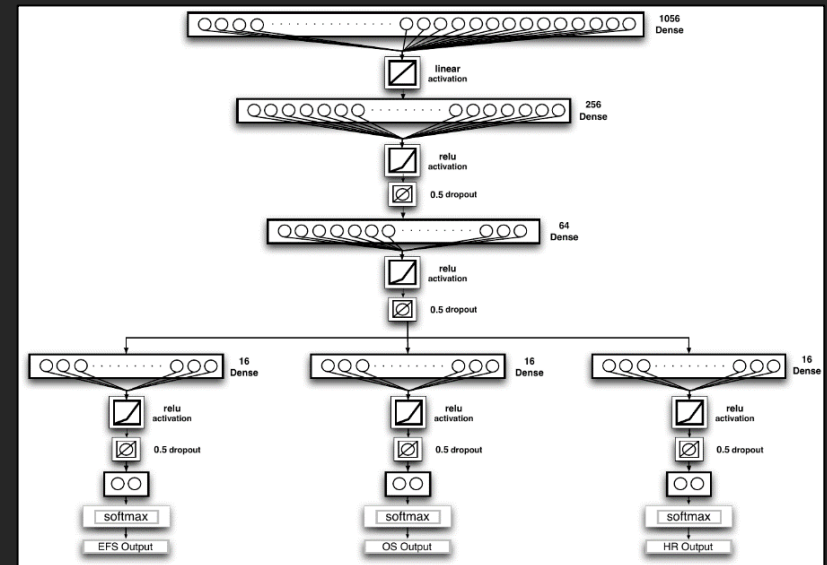
- **Deep Learning as accelerator of new ideas (e.g. embeddings, GANs)**

# MPBA: DEEP LEARNING FOR MASSIVE DATA



**Multi-Objective Deep Learning on Massive Sequencing Data**

**A. SEQC 2017: QC of massive NGS data for Precision Medicine**

**B. Food safety**

# Deep Learning & Omics

New Results    Posted May 28, 2017.

**Opportunities And Obstacles For Deep Learning In Biology And Medicine**

Travers Ching, Daniel S. Himmelstein, Brett K. Beaulieu-Jones, Alexandr A. Kalinin, Brian T. Do, Gregory P. Way, Enrico Ferrero, Paul-Michael Agapow, Wei Xie, Gail L. Rosen, Benjamin J. Lengerich, Johnny Israeli, Jack Lanchantin, Stephen Woloszynek, Anne E. Carpenter, Avanti Shrikumar, Jinbo Xu, Evan M. Cofer, David J. Harris, Dave DeCaprio, Yanjun Qi, Anshul Kundaje, Yifan Peng, Casey S. Greene

doi: https://doi.org/10.1101/142760

Deep learning-based methods now match or surpass the previous state of the art in a diverse array of tasks in patient and disease categorization, fundamental biological study, genomics, and treatment development. Returning to our central question: given this rapid progress, has deep learning transformed the study of human disease? Though the answer is highly

---

New Results    Posted March 8, 2017.

**Deep Learning based multi-omics integration robustly predicts survival in liver cancer**

Kumardeep Chaudhary, Olivier B. Poirion, Liangqun Lu, Lana Garmire

doi: https://doi.org/10.1101/114892

Auto-encoders + SVM

---

## CHANGING THE COURSE OF GENOMIC MEDICINE

Today, machine learning and experimental biology are advancing at an exponential pace. Deep Genomics is where these disciplines meet. Our systems predict the molecular effect of variation, opening a new and exciting path to discovery for disease diagnostics an...

LEARN MORE

DEEP LEARNING    GENOMICS    PRECISION MEDICINE

---

**Briefings in Bioinformatics**

Issues    Advance articles    Publish ▾    Purchase    Alerts    About ▾

**Deep learning in bioinformatics**

Seonwoo Min, Byunghan Lee, Sungroh Yoon

Brief Bioinform bbw068.    DOI: https://doi.org/10.1093/bib/bbw068
Published: 25 July 2016    Article history ▾

Although deep learning holds promise, it is not a silver bullet and cannot provide great results in ad hoc bioinformatics applications. There remain many potential challenges, including limited or imbalanced data, interpretation of deep learning results, and selection of an appropriate architecture and hyperparameters. Furthermore, to fully exploit the capabilities of

---

608    IEEE TRANSACTIONS ON NANOBIOSCIENCE, VOL. 14, NO. 6, SEPTEMBER 2015
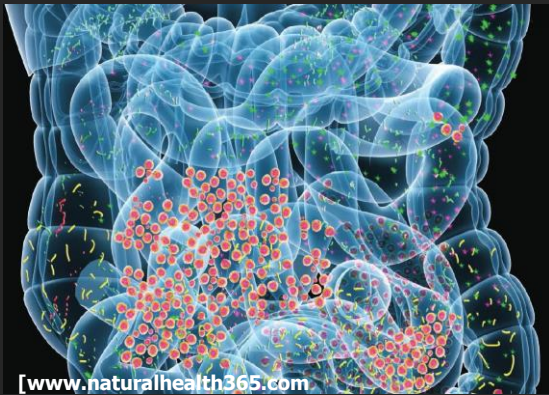
**Multi-Layer and Recursive Neural Networks for Metagenomic Classification**

Gregory Ditzler*, Member, IEEE, Robi Polikar, Senior Member, IEEE, and Gail Rosen, Senior Member, IEEE

The experiments discussed in the previous section demonstrated that: i) the deep learning approaches are not superior, at least on the data sets we evaluated, and ii) traditional MLPNNs are quite competitive with the RFCs, and in general perform better. However, none of the classifiers—deep or shallow—uniformly performs better than the RFCs across different experiments. The performance of the deep learning approaches may be improved upon with data sets that are much larger. It appears that—at least based on accuracy alone—the deep learning approaches may not be suitable for metagenomic applications. Accuracy, however, is not the only figure of merit.
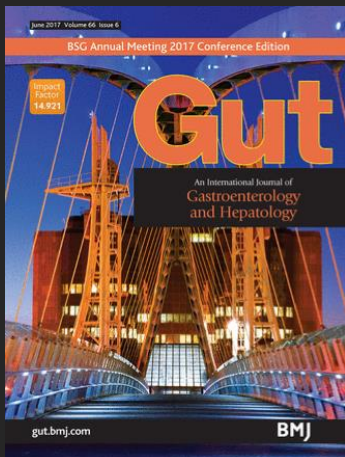
# Metagenomics for Gut Inflammatory Disease
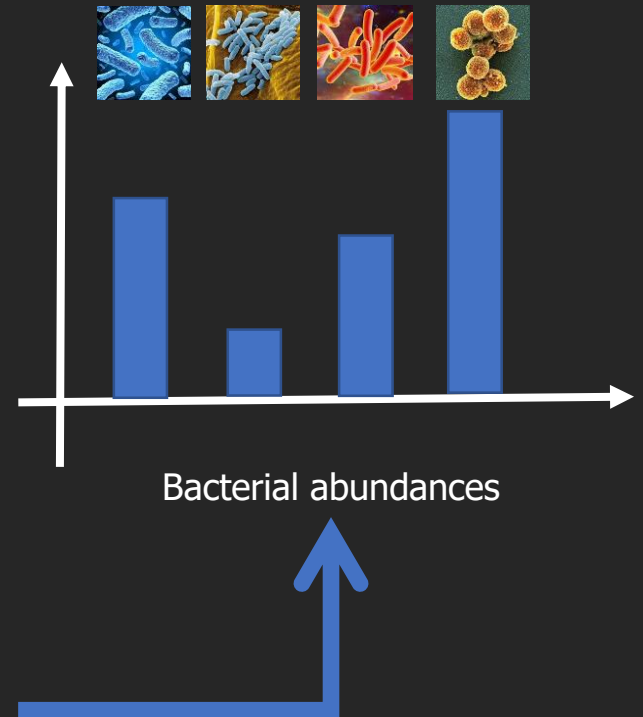


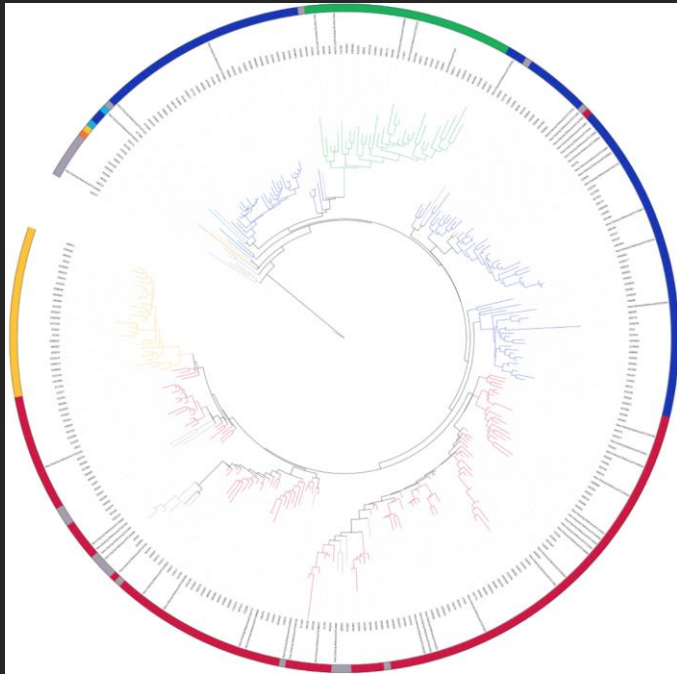Gut microbiome



Next Generation
Sequencing

Bioinformatics
pipeline



Bacterial abundances

Inflammatory
Bowel Disease,
with H Sokol

*C. Furlanello – MPBA Nov2017*

# A CNN architecture for Metagenomics

**Phylogenetic tree**

**Distance matrix**



**Convolutional Neural Networks**

[Perone, 2016]

# Defining a tree distance from "Patristic"

**Key ingredient:**
Concept nearest neighbours    $d($   ,  $)$

**Patristic distance:**

$\sum$(length of all branches) connecting two species

# Ph-CNN



Coordinates

Bacteria
Abundancies

**Phylo-Conv Layer
(4x16)**

**Phylo-Conv
(4x16)**

**MaxPooling
(2x1)**

**FC
(64)**

**Dropout
(0.25)**

**FC**

K-nearest
Phylogenetic
Neighbours

Conv +
SeLU

IMPLEMENTATION
- A new Keras layer
-    TensorFlow
-    GPU Cloud
-    ADAM optimizer

*C. Furlanello – MPBA Nov2017*

# GENERALIZATION TO OmicsCNN



Fioravanti et al 2017, BMC
Bioinformatics, in press

Jurman et al 2017 "Convolutional neural networks
for structured omics: OmicsCNN and the
OmicsConv layer

Machine Learning in Computational Biology
NIPS 2017 (9 Dec)

**arXiv:1710.05918**

# Deep Learning for Diagnosis and Prognosis of Pediatric Cancer (500 RNA-Seq)



Maggio et al 2017 " MULTIOBJECTIVE Deep Learning Approach for predictive classification in Neuroblastoma"
ML4H: Machine Learning for Health Workshop NIPS 2017 (Dec 8, 2017)

# PRECISION MEDICINE

**A unifying framework**

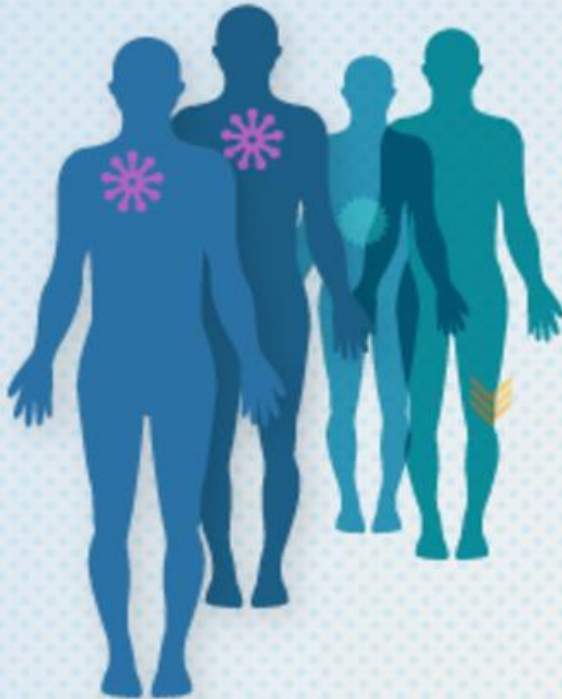• *"Treating the **right patient** for the **right disease** at the **right time** with the **right amount of (the right) drug"***
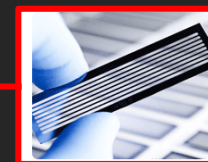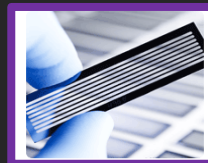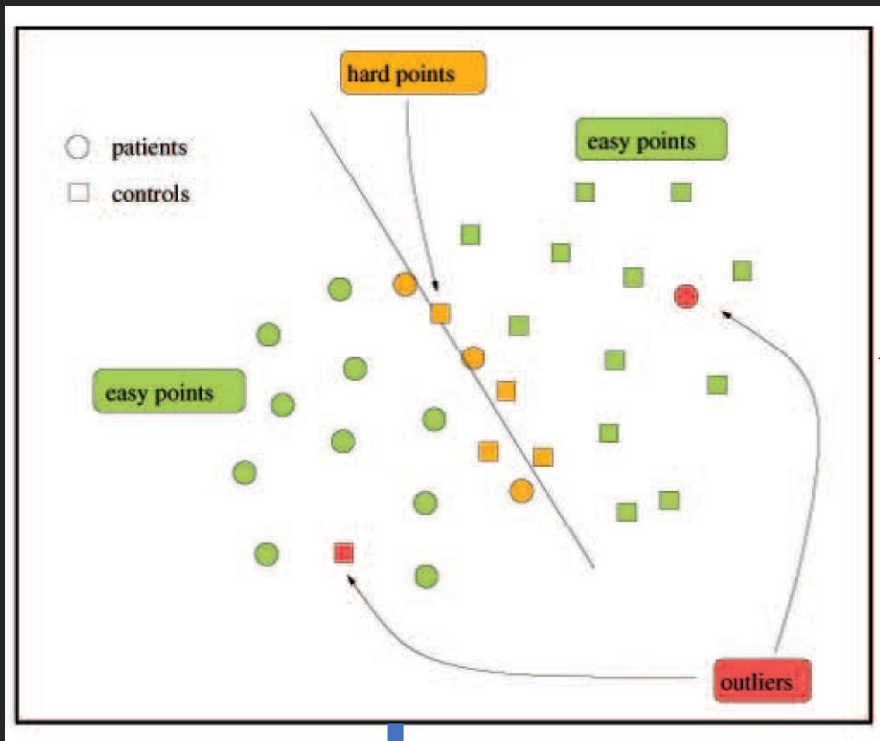
*In Diagnostics, treatment and prevention,* we shall systematically include the **individual variability** of **genes, environment, life style**
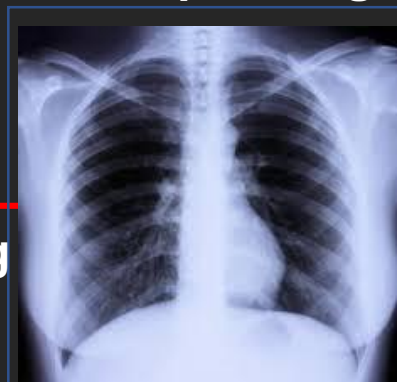
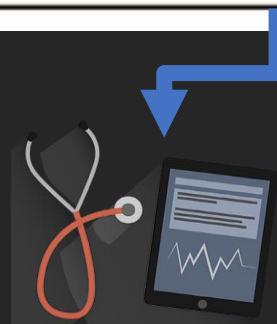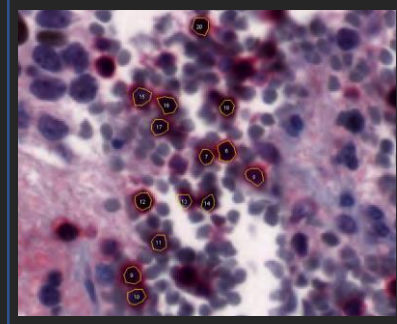The Precision Medicine Initiative (USA 2015) is the vision that redefines healthcare and Pharma R&D

# MACHINE LEARNING IN BIOMEDICINE

**Next Generation Sequencing**

**Bioimaging**

**Outcome
(complex phenotype)**

**LAB**

# CHALLENGE

**Genomic Medicine R&D is a patchwork** in terms of data, technologies and models hard to integrate. We are studying a Deep Learning solution to identify a common share space for health trajectories



Multi-omics
(Biomarkers and Networks)

Diagnostic Bioimaging
(Pathology and
RADIOLOGY)

**Multiple Endpoints**
(Diagnostics, prognostic, actionable)

Modulated by
**Complex Clinical Phenotypes**

Phenotypes
(Lab, EHR, ....)

Natural integration in
the latent space

Deep Learning to recover data still unconnected, of highest value for reconstruct the entry gates to pathology, or to intercept their similarity, or transferable therapy between diseases.

# FROM HEALTH TO ENVIRONMENT, AND BACK

## Spatio-temporal NOWCASTING (Conv-LSTM)

# NOWCASTING

Rain & lightning nowcasting  (5' - 75' )
Short-time radar prediction: target for **deep learning**

# Deep Learning in AgriTech



Expert estimating damage

- **Time-consuming methods based on limited data**
- **High variance**



Faster-RCNN ~54M pars

1. **Yield quantification**
2. **Assessment of damage and risk**
3. **Quality Control**



Tested July 2017, Val di Non, Trentino

# Deep Learning as a Service

# Deep Learning for retail

**Embeddings**
(categorical)

**Sales**
(Numeric)

input → output

**Target:** daily sellout forecast (60 days)

**What**: Deep Learning pipeline for POS sales

**First large scale application:** pharmaceutical retail (UNIFARM),

- **315K unique products, 458 POS** over **1642 days** (~4 years),

- Total sales of **€ 2.4 Billions**

- The module is **accurate** (improves over XGBoost): **16% RMSPE**

eit Digital

Partners: FBK, REPLY, Tim, BT, DFKI Saarbrucken …

# DL User profiling (fintech security)

**BIOLOGICAL BIOMETRICS**

Your iris & retina   Your fingerprints   Your DNA

**BEHAVIORAL BIOMETRICS**

**Keystroke dynamics: how a user types**

- **Monitor keyboard inputs (KHz),**

- **Process typing features**

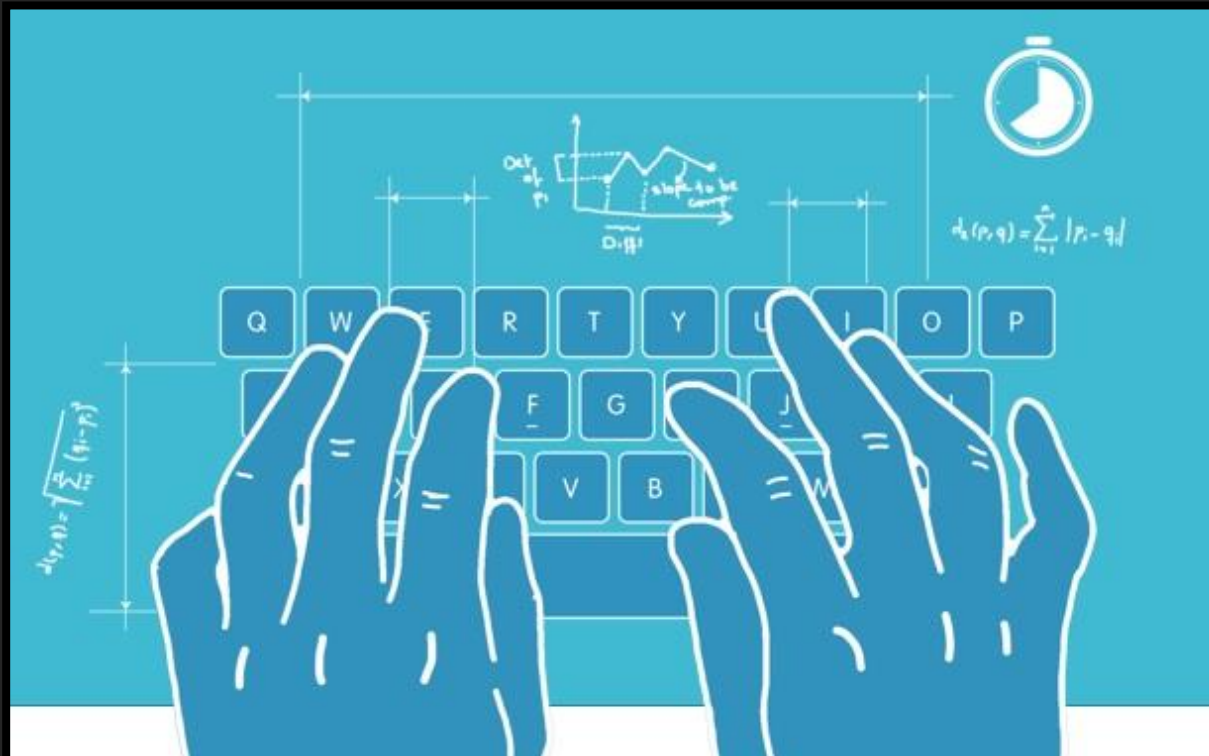- **Defines a pattern for future comparison**

Your writing & signature   Your voice & speech patterns   **Your typing speed & patterns**

DWELL TIME
The time between pressure and release of a key

GAP TIME
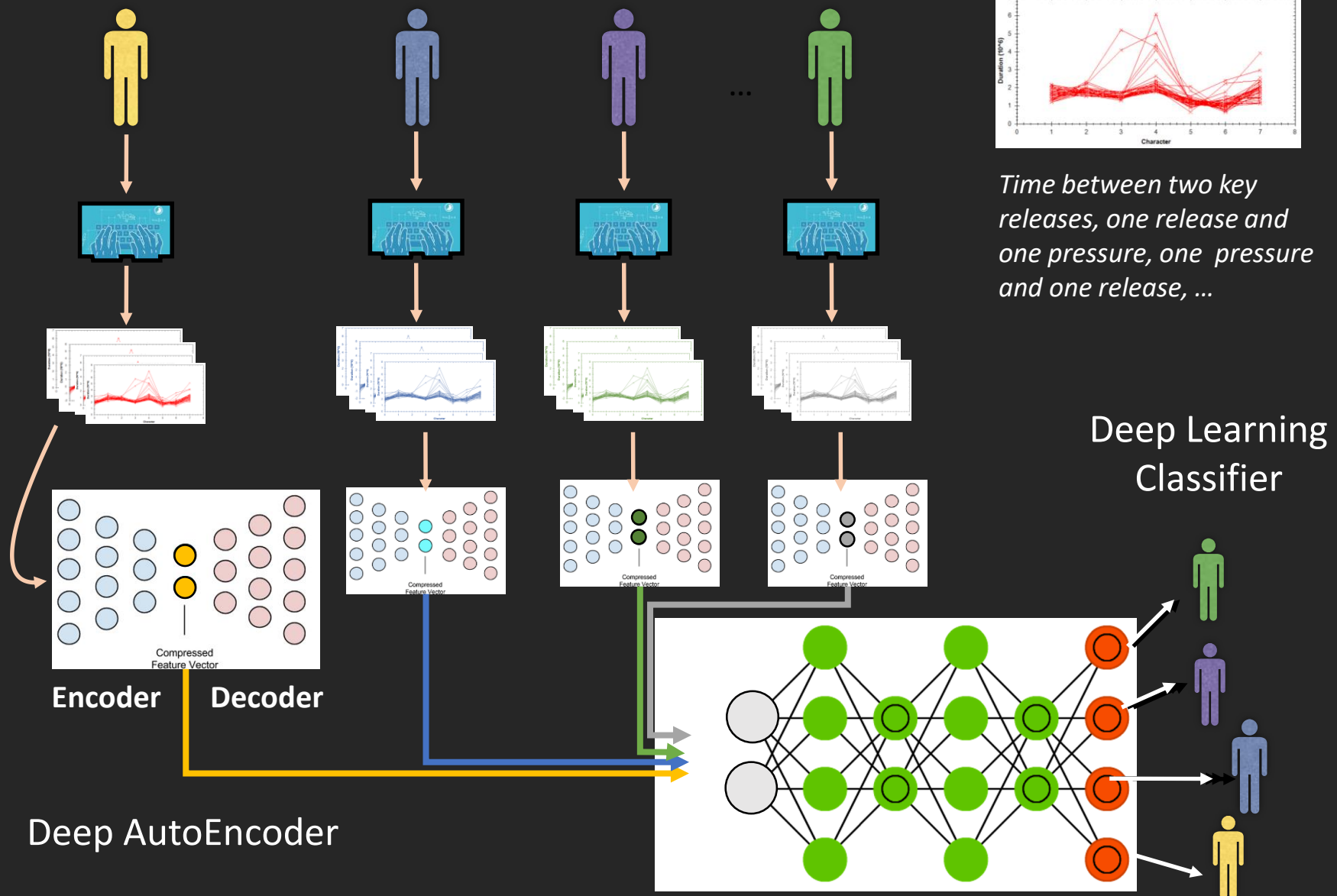The lapse of time between striking one key and the next

Your typing speed & patterns

Identifying an individual based on her way of typing on a physical or virtual keyboard

eit Digital

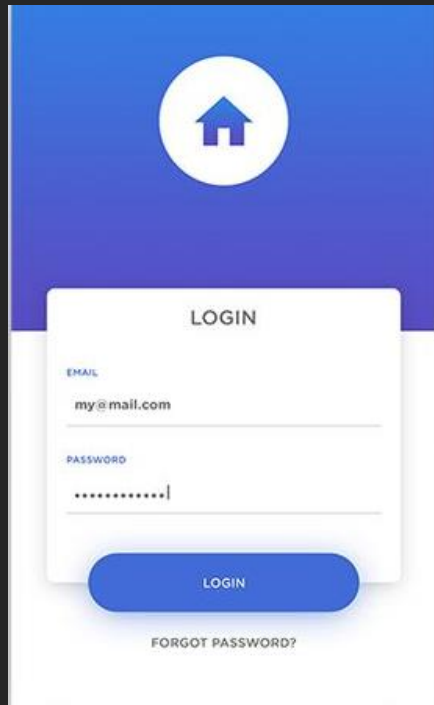Partners: FBK, IMDEA, REPLY COMMUNICATION VALLEY

# Individual DeepKS Learner

*Time between two key releases, one release and one pressure, one pressure and one release, ...*

**Deep Learning Classifier**

**Encoder** | **Decoder**

Deep AutoEncoder

# Deep Learning as a Service



API Engine

Feature extractor

DeepKS Model

{json}

Eve
python-eve.org

github.com/
spotify/luigi

Raw data, features,
predictions

MongoDB

*C. Furlanello – MPBA Nov2017*

eit Digital

# Execution

**R&D competition driven by a new vision**

**Innovation in DL**

**"run as a venture" within collaborations**

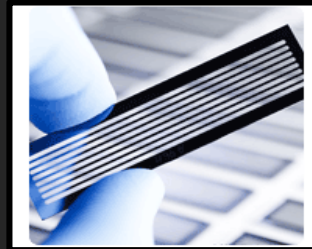**Commit to grow a new generation** of interdisciplinary researchers with strong entrepreneurship

**Access to GPU resources, datasets, & "open science setups"**

*C. Furlanello – MPBA Nov2017*
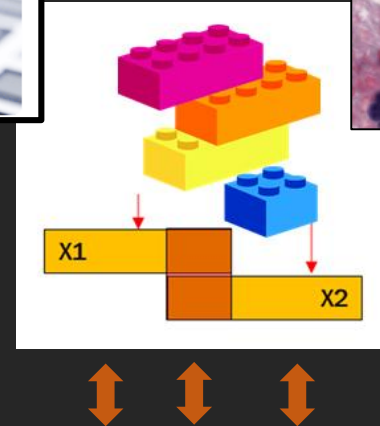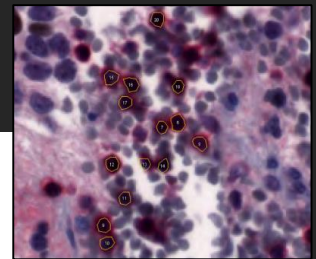
Resources for disruptive innovation

# Unexpected Synergies

- **IOT for Smart Digital Industry**
- **Precision Medicine** for biotech/healthcare research
- **Environment:** big data and AI for life and food

## Acceleration of Deep Learning technology on Agritech



Enabling precision medicine on multimodal data

Bayer and Monsanto – 2017

**WebValley is the FBK summer school for data science and interdisciplinary research: close to 350 students from around the world (17-19y old) have attended the WebValley camps since its first edition in 2001.**

**In 2016 and 2017, the team developed a new Deep Learning solution for fruit quality control based on portable spectrometry and low cost images**

**Agritech as an accelerator of Precision Medicine: Deep Learning, cloud infrastructure (MS Azure), local GPU boxes, blockchain**

# START FAST, START EARLY