# From the Higgs to Huntington's: methods for learning from data

**UCL HEP seminar**
**24-05-19**

**Peter Wijeratne**
**MRC Skills Development Fellow**

**UCL CMIC**

Neil Oxtoby
Alexandra Young
Arman Eshaghi
Leon Aksman
Maura Bellio
Nonie Alexander

**UCL HDC**

Sarah Tabrizi
Rachael Scahill
Sarah Gregory
Eileanoir Johnson
Ed Wild
Lauren Byrne

**CHDI**

Cristina Sampaio
Amrita Mohan
John Warner
Dorian Pustina
Alexandra Shechtel

And all the participants of the PREDICT, TRACK and IMAGE-HD studies.

Interested in extracting hidden information from observed data

$\rightarrow$ Bayesian methods
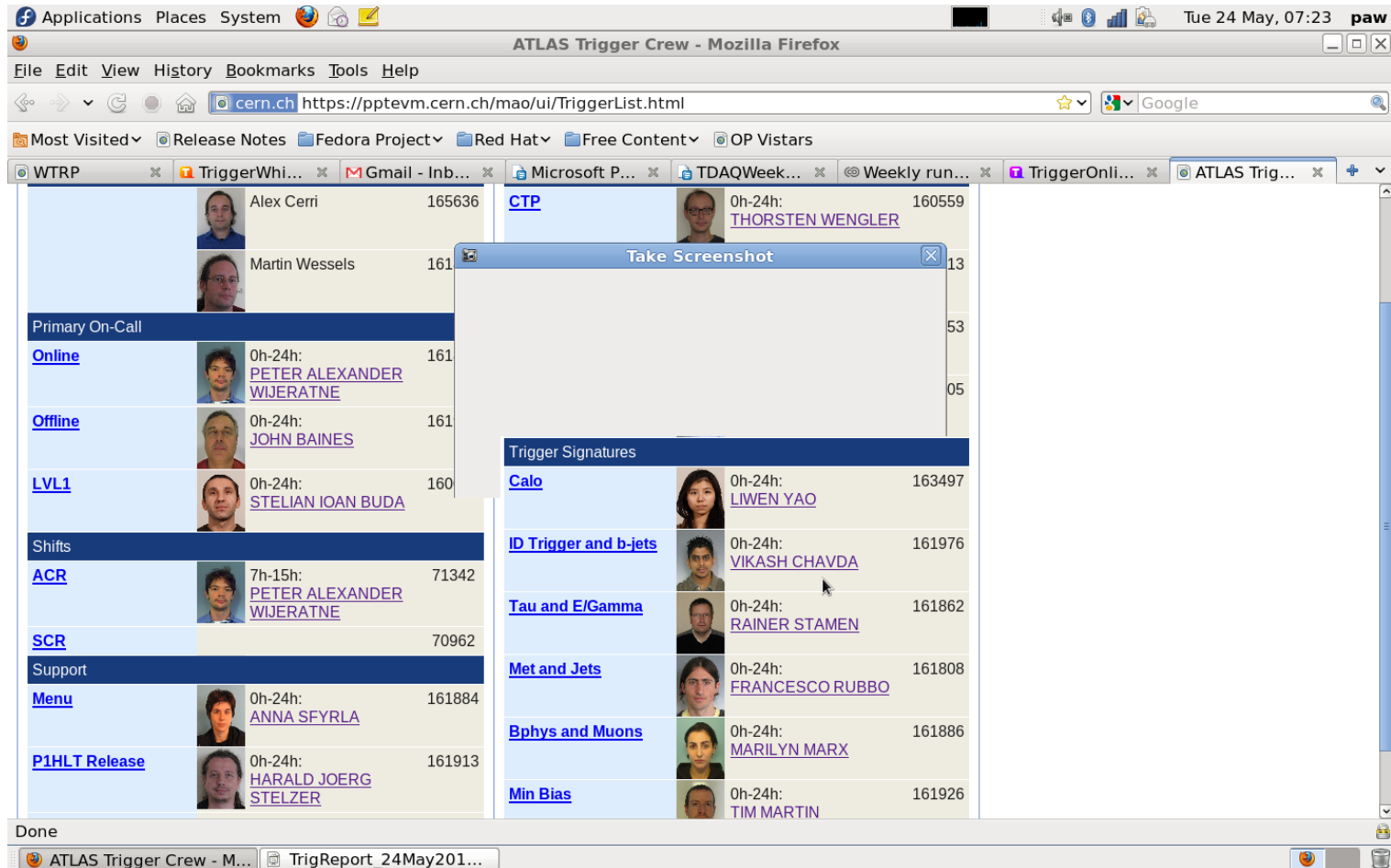
Two main schools of thought

Hypothesis-driven (informative priors)

Unfolding / inverse problems – e.g. image reconstruction

Data-driven (non-informative priors)

Latent variable inference – e.g. disease progression modelling

Physics favours the former, biology the latter

For some reason, they let me near the detector

Nature ← Unfolding ← Measurement

$$n(C_i^{data}) = \frac{1}{\epsilon_i} \sum_j P(T_i^{MC} \mid R_j^{MC}) n(R_j^{data})$$

• Real data are dependent on the detector used to measure them

• Bring data back to their natural state by applying hypothesis-driven corrections derived from simulation

→ "Unfolding the cause"

- Energy density (min bias + UE) was not modelled correctly in forward direction

    - Problem would only increase with luminosity

- We iteratively unfolded the data to compare directly with various models

- Tuned MC generators to data

4

I saw this one day in 2013



I wanted to use physics to fight cancer

I asked about for potential opportunities (thanks Simon)

I got lucky and a postdoc came up at the Centre for Medical Image Computing on jobs.ac.uk

Maths, physics and engineering scientists at the interface of basic and biomedical sciences



CMIC

Great Ormond Street Hospital
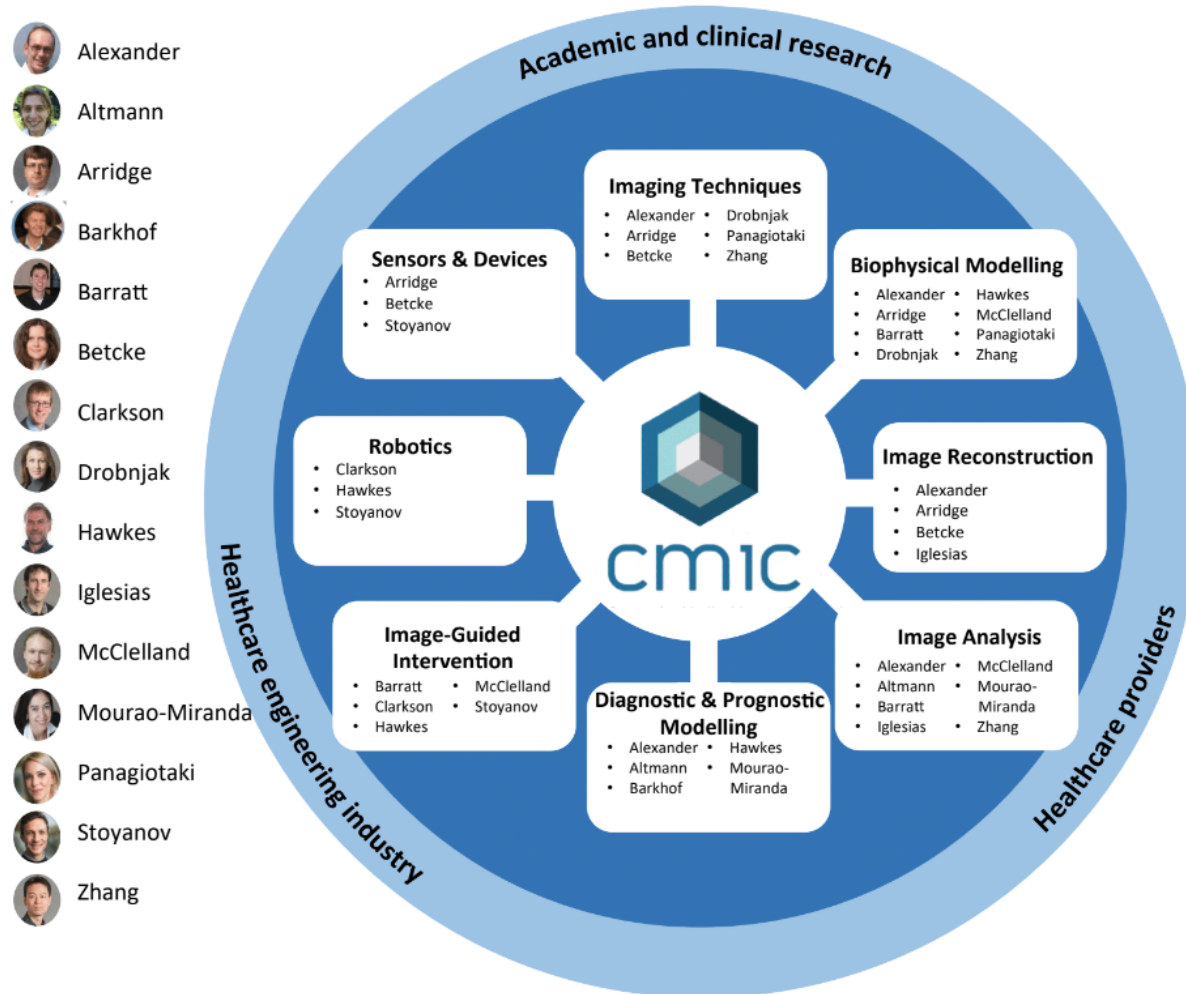
University College London Hospital
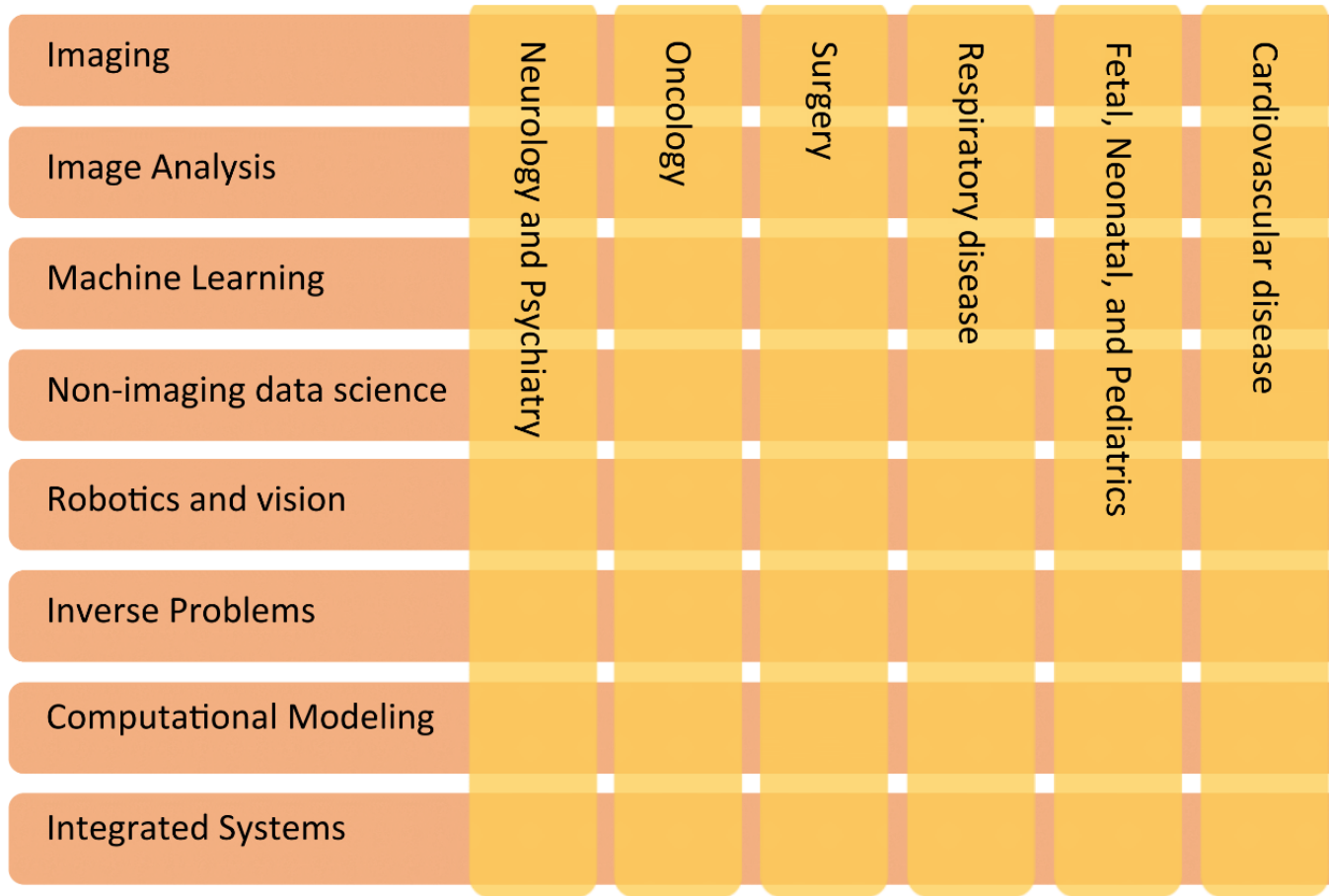
Moorfield's Eye Hospital

Royal Free Hospital

Royal National Orthopaedic Hospital
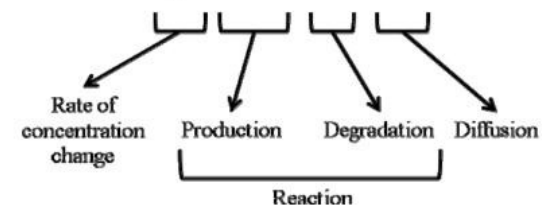
**The Chemical Basis of Morphogenesis**

A. M. Turing

*Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, Vol. 237, No. 641. (Aug. 14, 1952), pp. 37-72.

$$\frac{\partial u}{\partial t} = F(u,v) - d_u v + D_u \Delta u$$

$$\frac{\partial v}{\partial t} = G(u,v) - d_v v + D_v \Delta v$$

Rate of concentration change    Production    Degradation    Diffusion
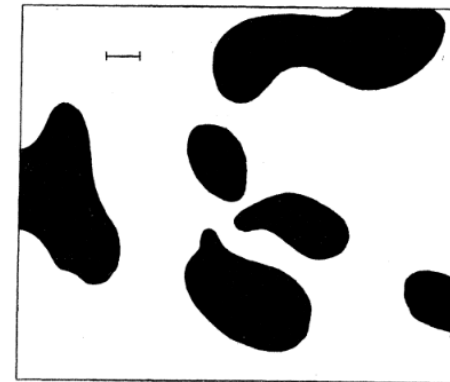
Reaction

Computational Modeling

FIGURE 2. An example of a 'dappled' pattern as resulting from a type (a) morphogen system. A marker of unit length is shown. See text, §9, 11.
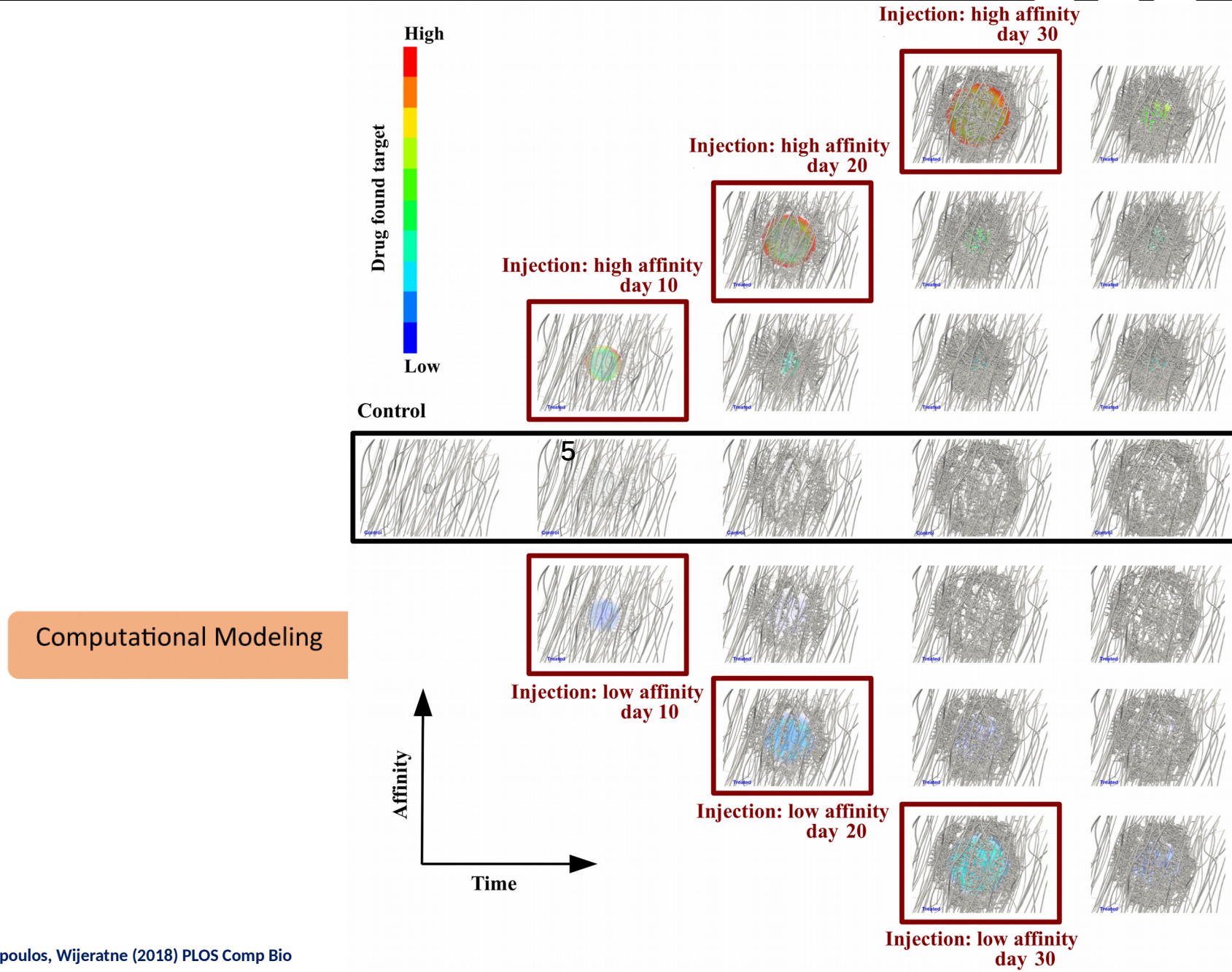
Injection: high affinity day 30

Injection: high affinity day 20

Injection: high affinity day 10

Drug found target — High / Low

Control

5

Computational Modeling

Affinity

Time

Injection: low affinity day 10

Injection: low affinity day 20

Injection: low affinity day 30

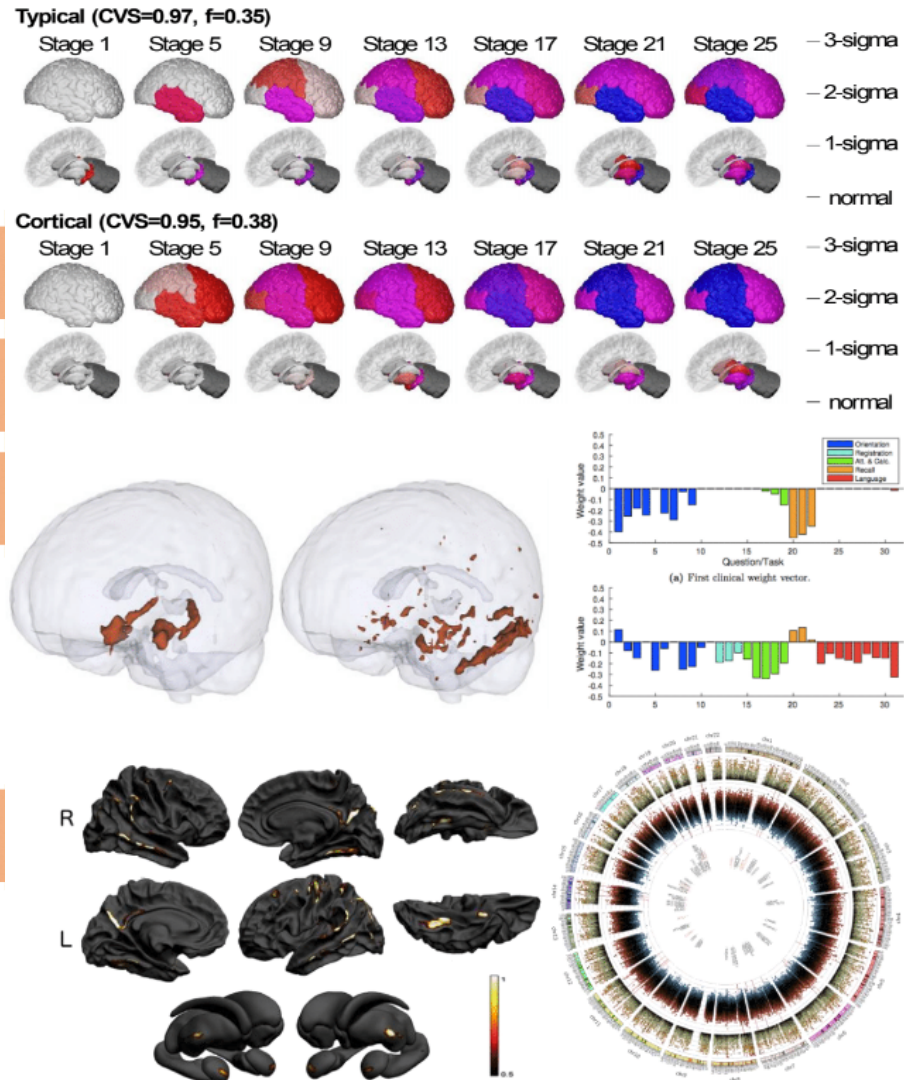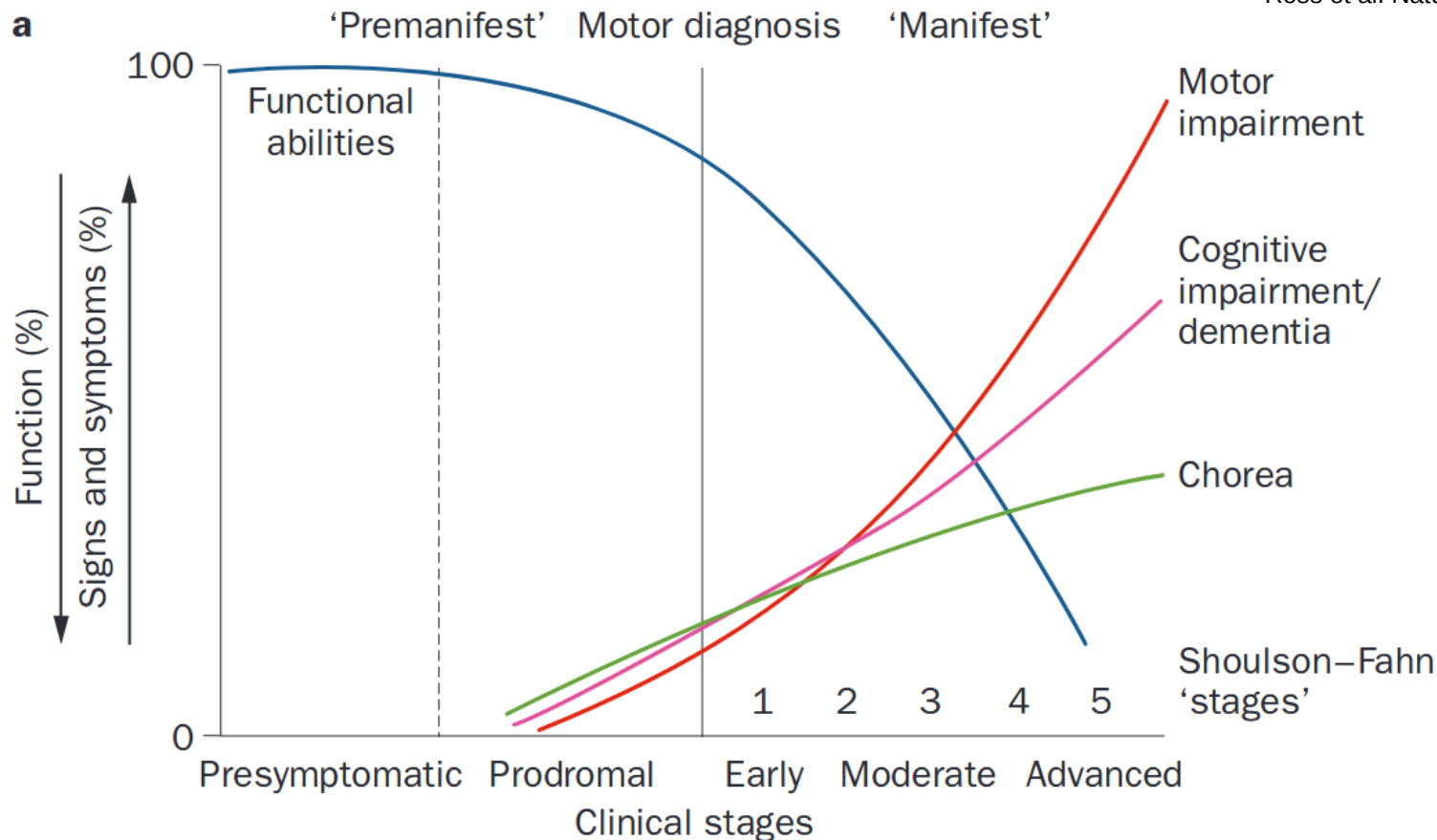Vavourakis, Stylianopoulos, Wijeratne (2018) PLOS Comp Bio

Image Analysis
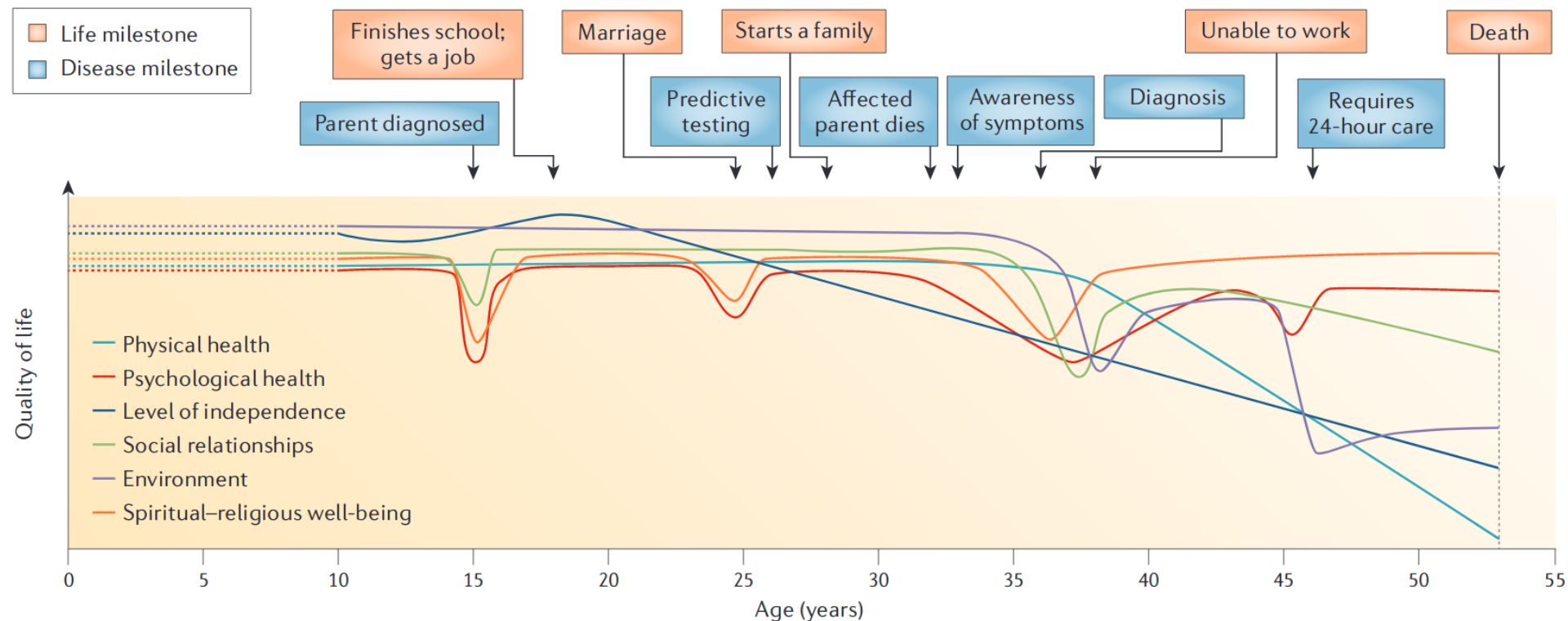
Machine Learning

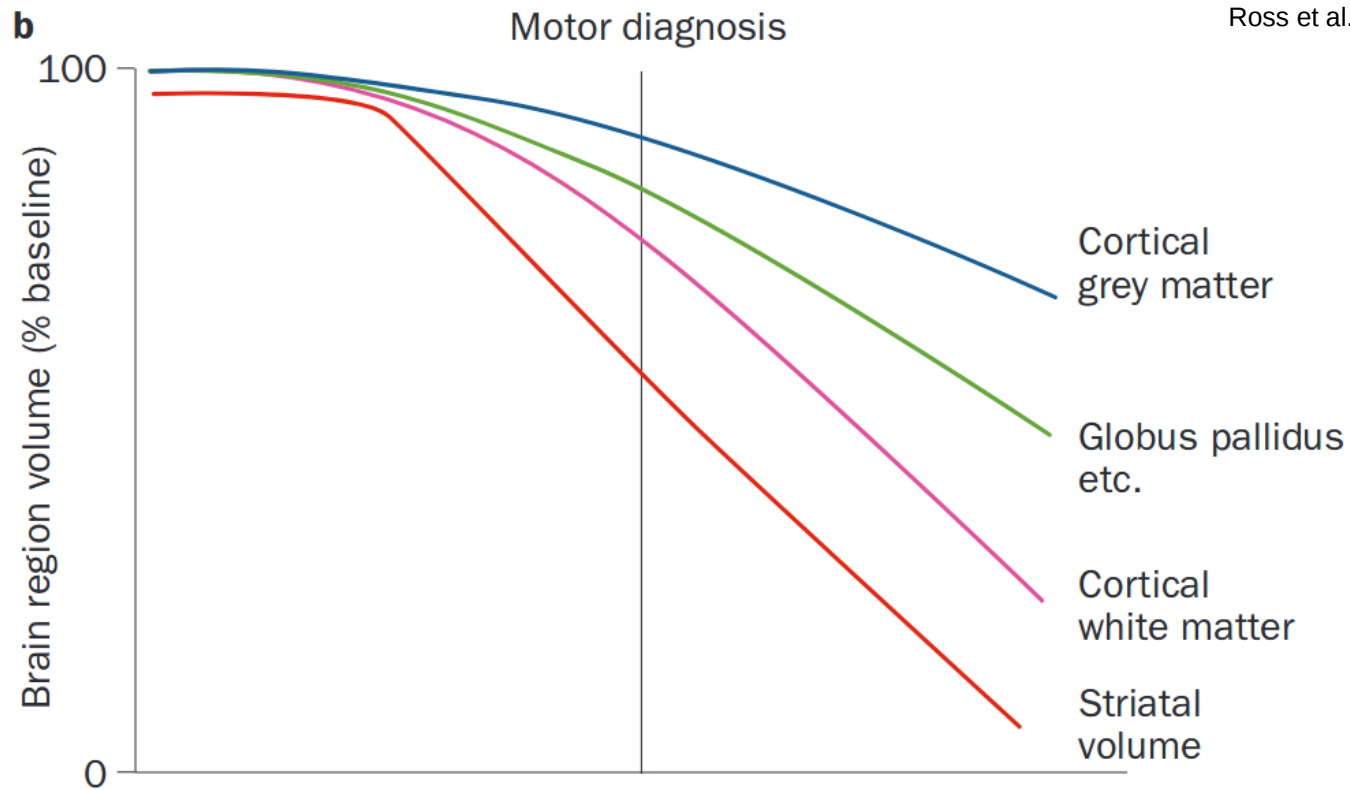Non-imaging data science

Computational Modeling

Slowly progressive, hereditary brain disease that causes changes in movement, thinking and behaviour
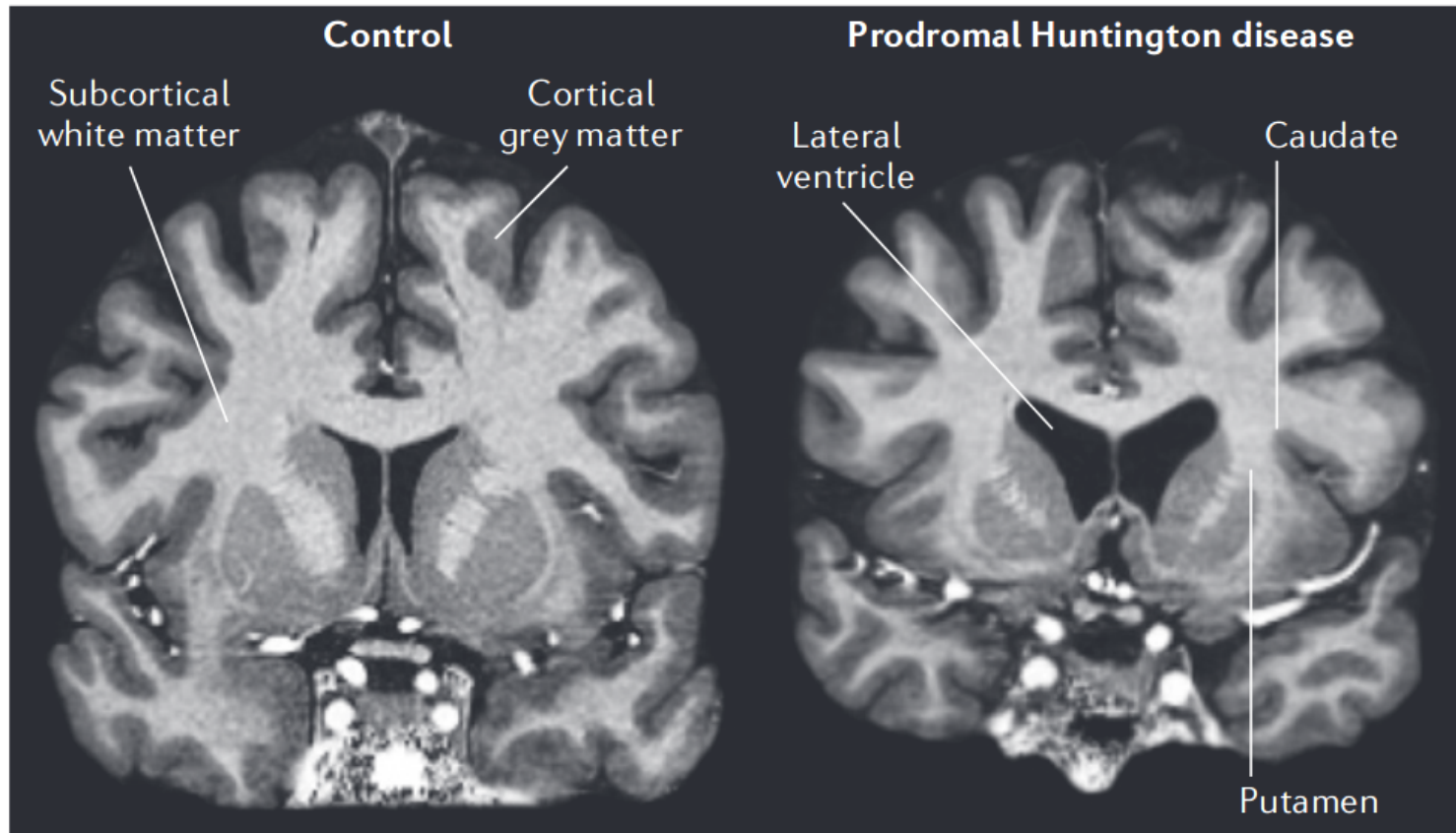
Autosomal dominant inheritance – 50% chance, everyone with gene will get HD

Diagnosis made at onset of movement disorder, typically with chorea and impaired voluntary movement
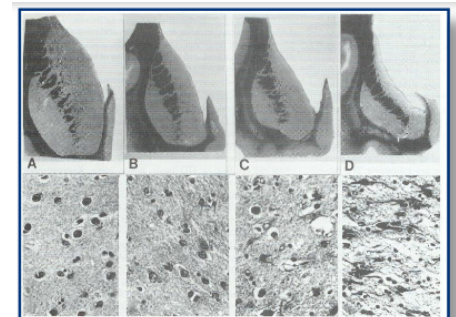
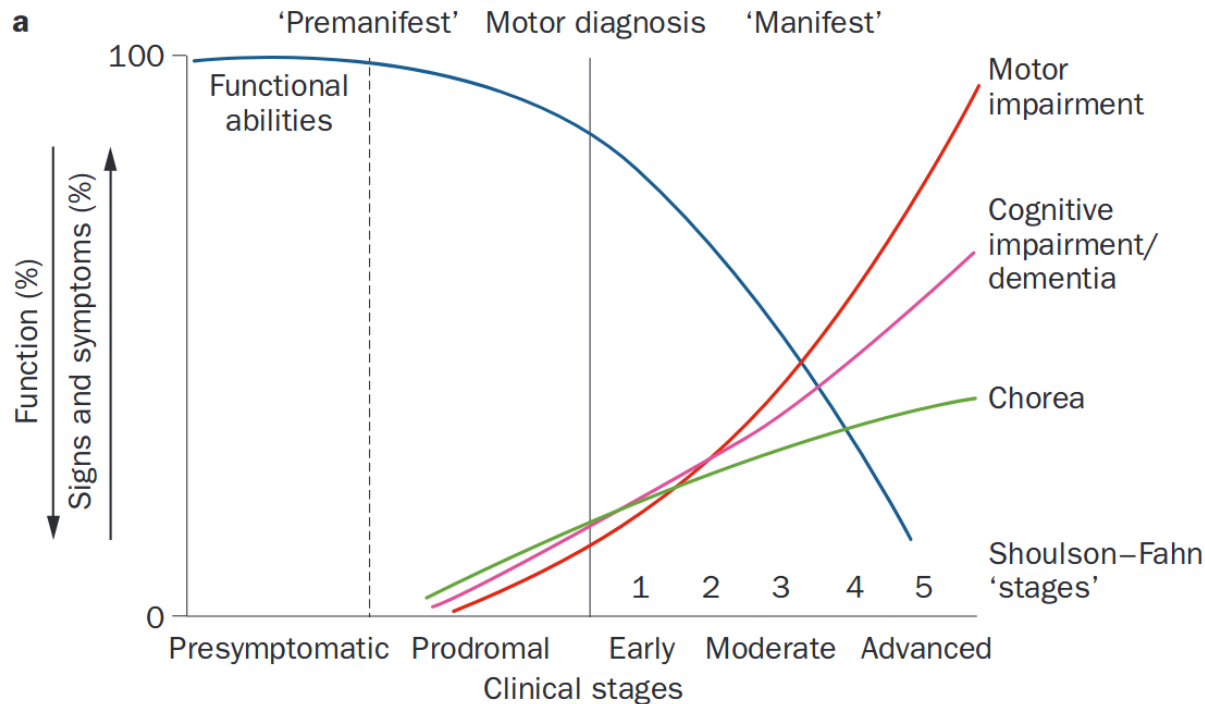Brain changes in HD – specific regions of the brain are atrophied

MRI provides spatial intensity measurements that depend on tissue properties

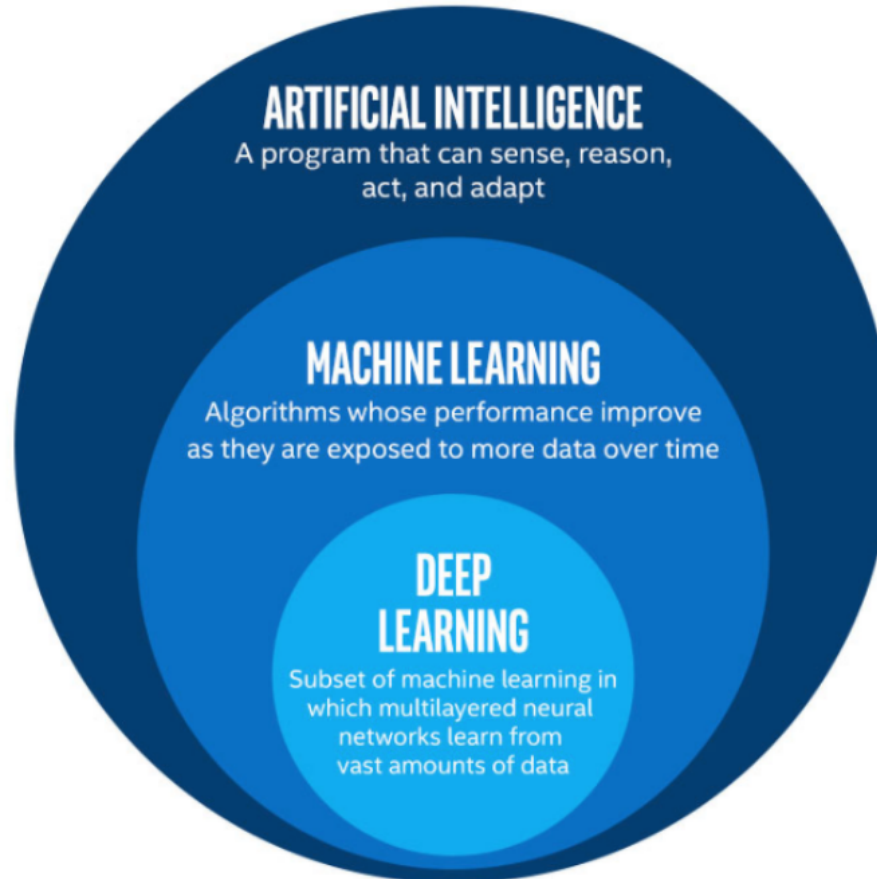Observed changes reflected by microscopy (histology)

15

Can we estimate where a patient is along their disease path?



Patient stage is a latent variable – it generates the observed measurements, but is not measured directly (unlike in physics events, where we know time)
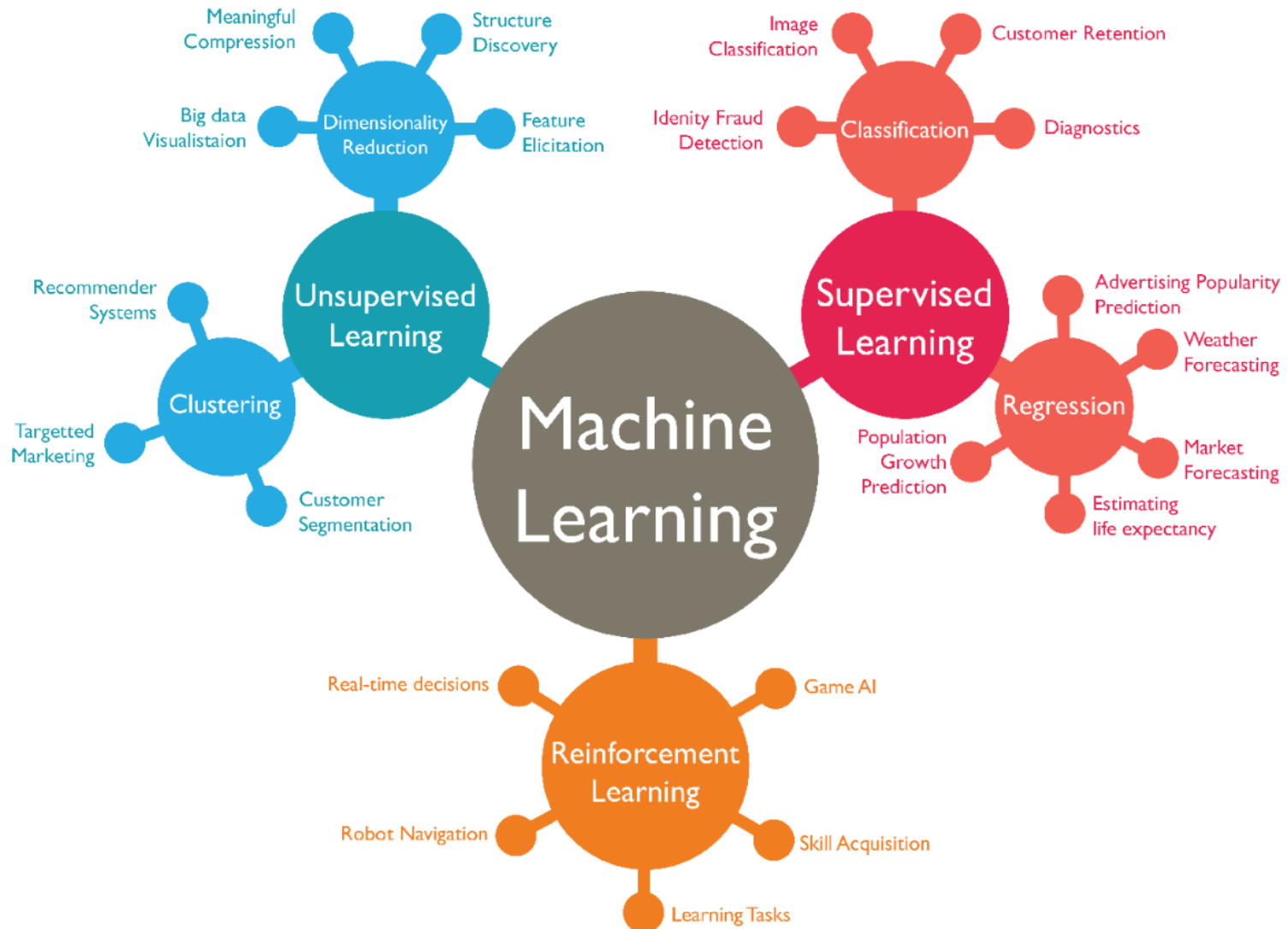
→ Infer using machine learning methods

http://www.learnwebskill.com/technology



Can think of machine learning as "data-driven AI"

Deep learning learns its own feature space
+ improved performance over standard ML methods
- difficulty in interpretability

What machine learning does well

1. Model-free identification of trends and patterns

2. Improves with data availability

3. Requires minimal (or no) human intervention
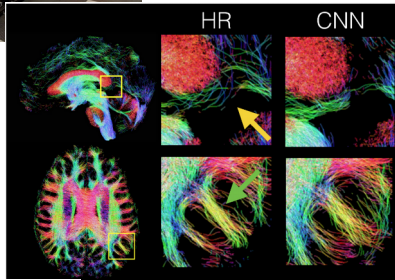
What machine learning doesn't do well

1. Causal mechanisms

2. Data intensive

3. Interpretability

We want to diagnose and prognose patients – don't really need to understand mechanisms
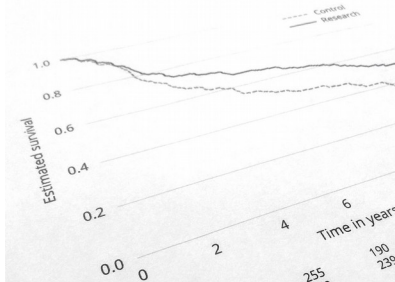
**Basic sciences**

**Clinical sciences**

**Cluster computing**

**Imaging + machine learning**

UCL EPSRC CDT in **Medical Imaging**

**Statistical methods**

20

Biomarker: any biological measurement that tracks disease progression

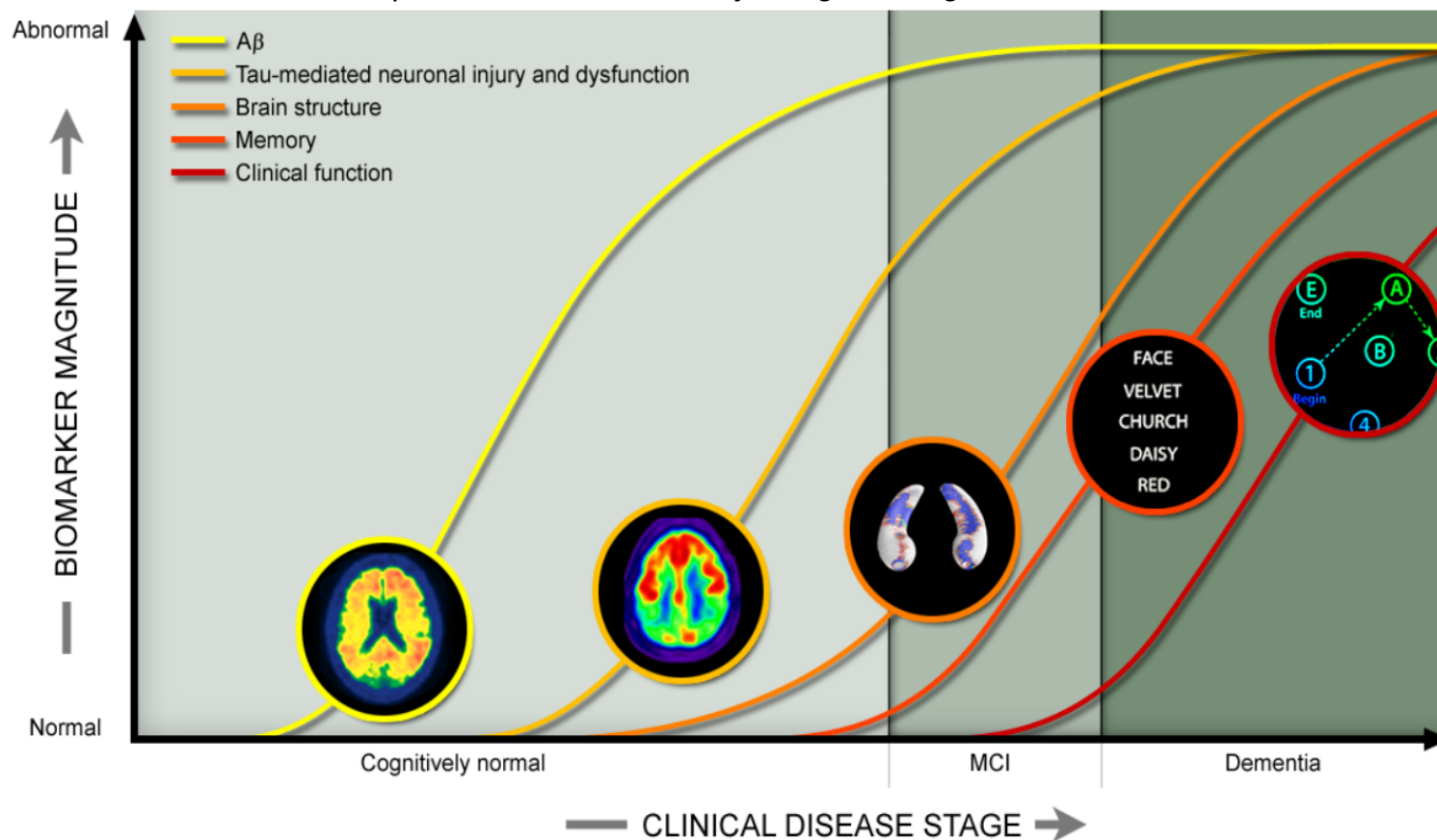Event: transition of a biomarker from a normal to abnormal state (Markovian)

Sequence: order of events over sample of interest

Cross-sectional: data from a single time-point

- Construct a picture of how disease plays out over time

- Express in terms of symptoms, pathologies and biomarkers

- Reconstruction must exploit cross-sectional data, where possible

http://adni.loni.usc.edu/study-design/#background-container

A picture of how components of a disease progresses over time

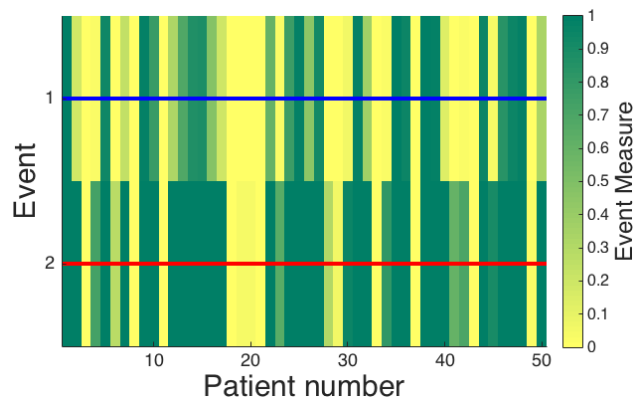Disease progression models learn patterns of disease-related changes from data



Machine learning

Patient data

Disease progression model

- Can use models to infer temporal ordering of changes

- Can also stage and stratify patients → clinical trial design

EBM estimates ordering of **binary events** from data – normal or abnormal

Data can be cross-sectional and any combination of types (imaging, clinical, genetic...)



$E_2$ $E_1$

Simple example: 2 event measures

More patients have greater abnormality in Event 2 than Event 1

→ Event 2 **measurably abnormal** before Event 1

More formally: EBM is a generative model of observed data from unknown sequence

data    uniform prior

$$P(X|S)=\prod_{j=1}^{N}\left[\sum_{k=0}^{Z}\left(P(k)\prod_{i=1}^{k}P(x_{ij}|E_i)\prod_{i=k+1}^{Z}P(x_{ij}|\neg E_i)\right)\right]$$

sequence                    prob.        prob.
                          Abnormal       Normal

- The EBM needs likelihood distributions for normal and abnormal subjects

→ Learn directly from data

▲UCL

Prince, SJD. Cambridge University Press. 2012

1. Mixture model fitting
– Expectation Maximisation



wikipedia.org/wiki/gradient_descent

2. Latent variable (sequence) fitting
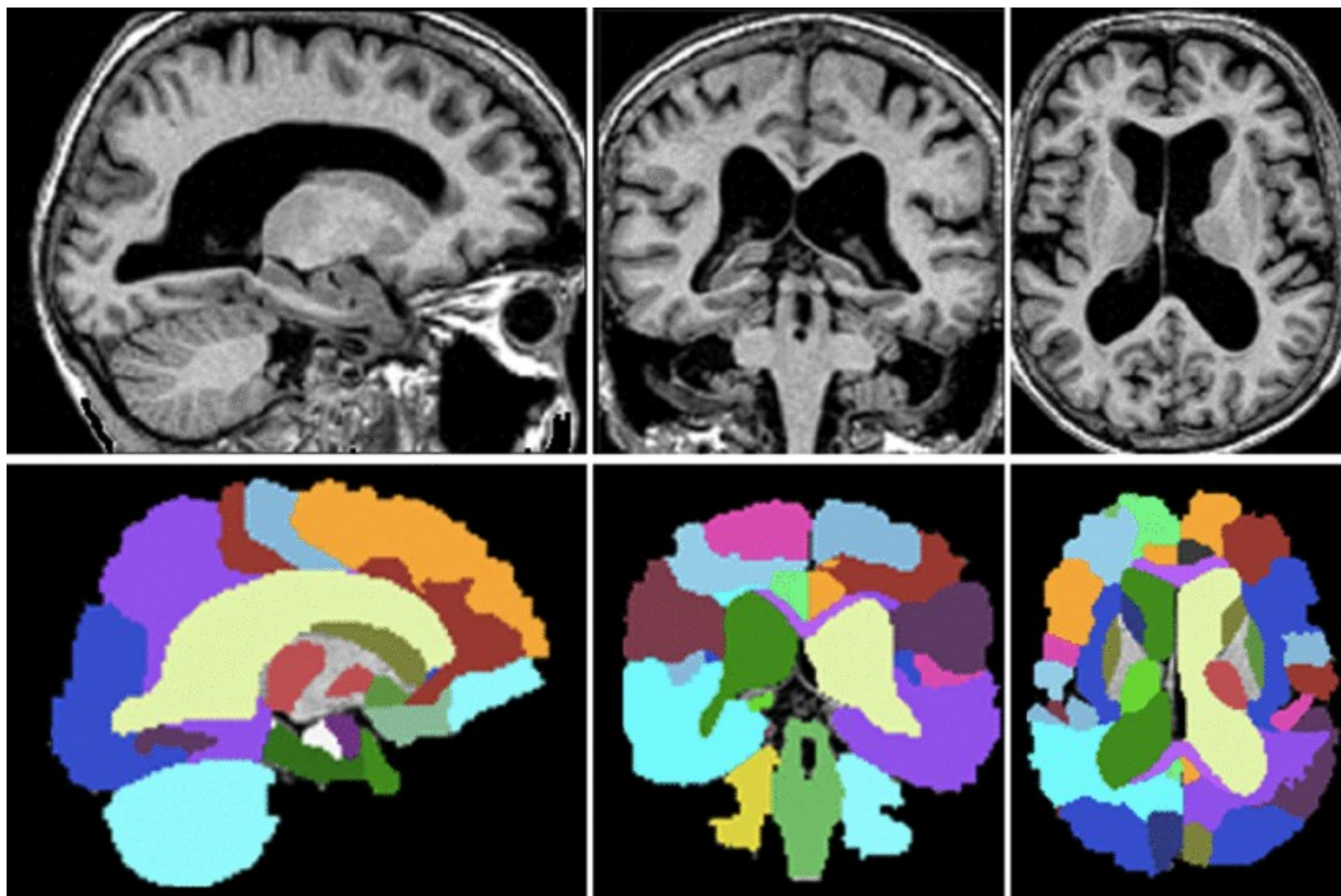– Gradient Ascent



3. Uncertainty estimation
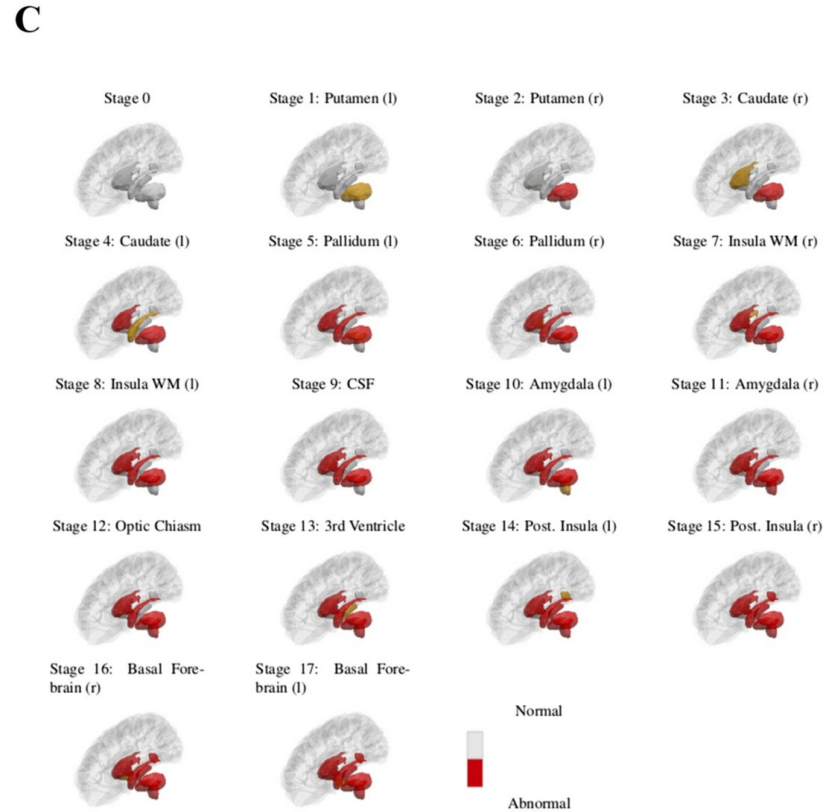– Markov Chain Monte Carlo
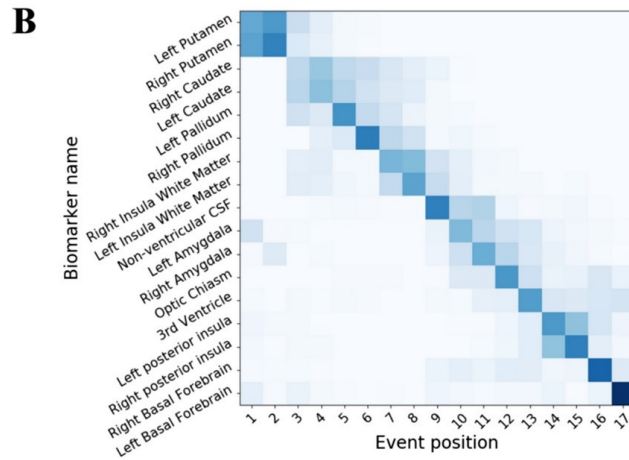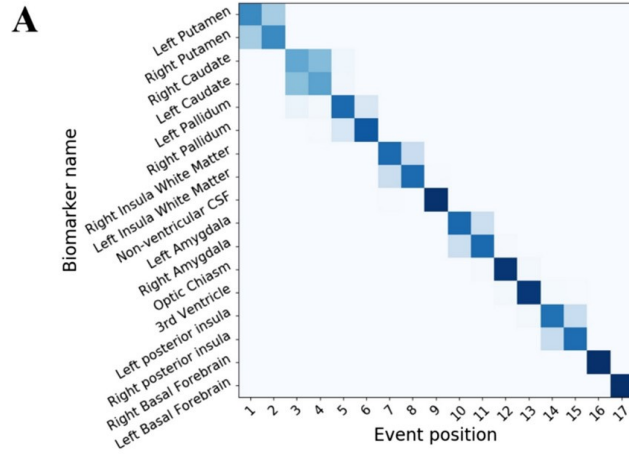


$$a = p(X \mid S')/p(X \mid S_t)$$

28

1. Build model on TRACK-HD

2. Cross-validate using PREDICT-HD and IMAGE-HD

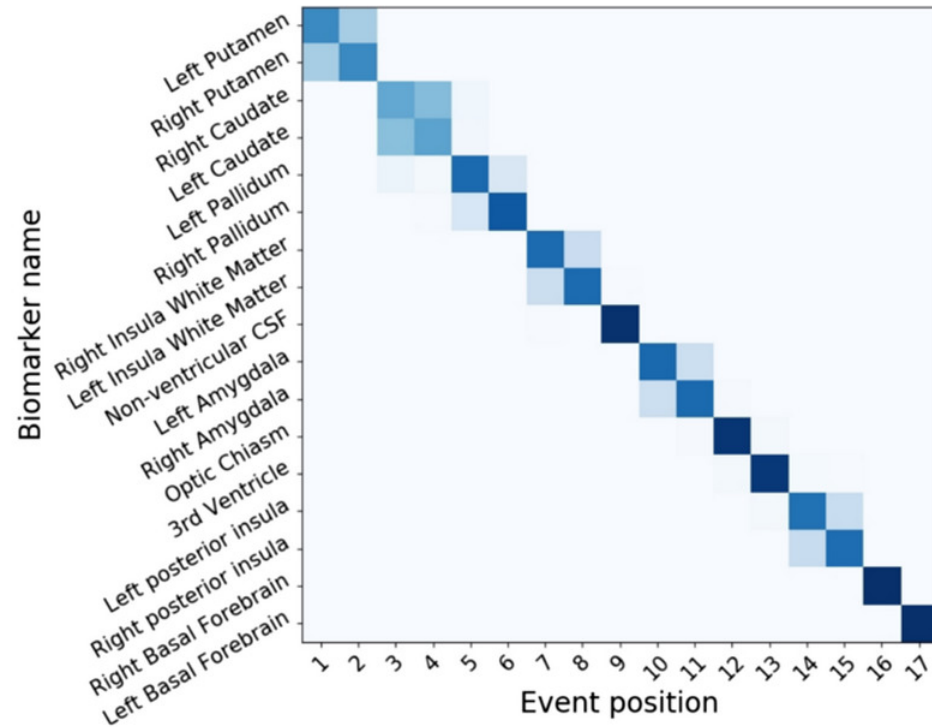3. Test predictive utility using TRACK-ON and PREDICT-HD

Extract regional brain volumes using Geodesic Information Flows*

→ Reduces inter-subject variability by using spatially variant graphs to connect morphologically similar subjects
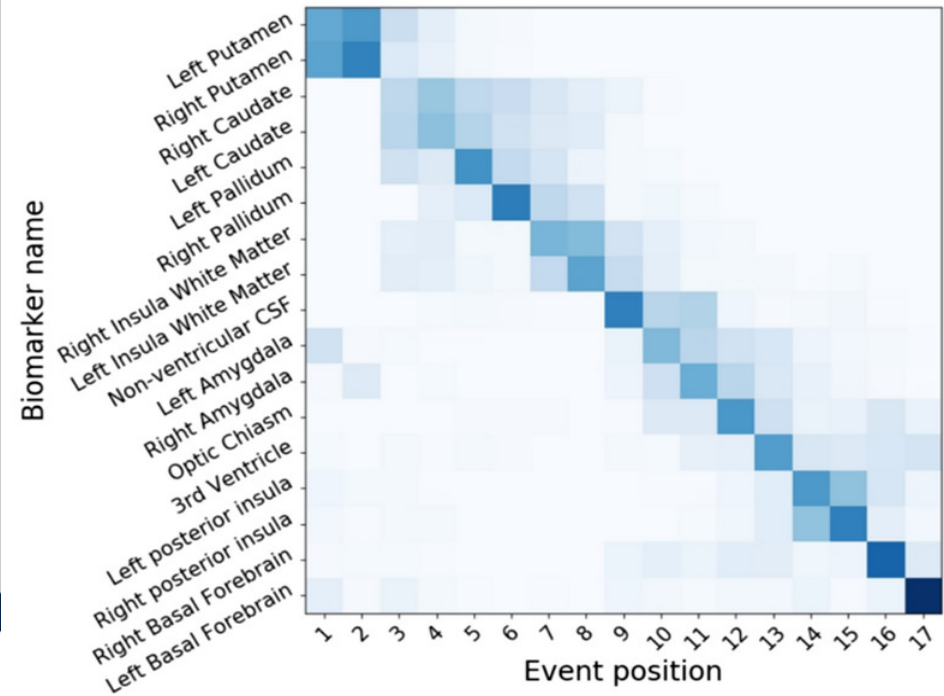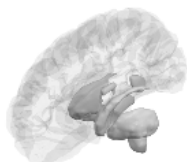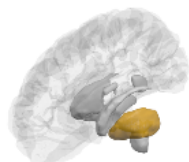
Direct model fit

Bootstrapped model fit

· Dark diagonal components indicate strong event ordering
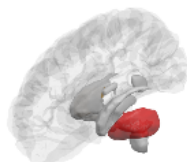
· Lighter indicate possible event permutations

32

Stage 0

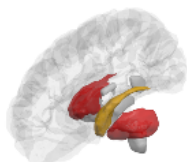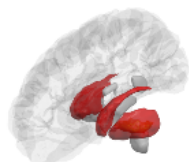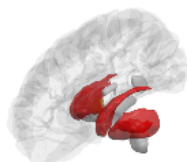Stage 1: Putamen (l)

Stage 2: Putamen (r)

Stage 3: Caudate (r)

Stage 4: Caudate (l)

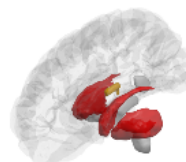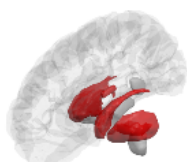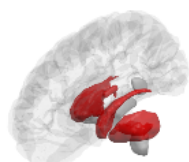Stage 5: Pallidum (l)
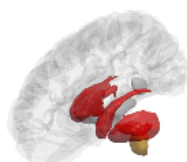
Stage 6: Pallidum (r)
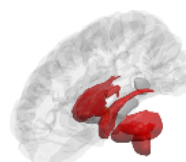
Stage 7: Insula WM (r)
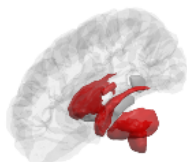
Stage 8: Insula WM (l)

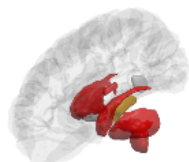Stage 9: CSF

Stage 10: Amygdala (l)
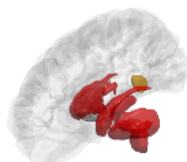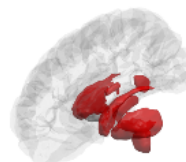
Stage 11: Amygdala (r)
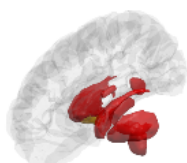
Stage 12: Optic Chiasm

Stage 13: 3rd Ventricle

Stage 14: Post. Insula (l)

Stage 15: Post. Insula (r)

Stage 16: Basal Forebrain (r)

Stage 17: Basal Forebrain (l)

Normal

Abnormal

Central

HD progression

Peripheral

SCIENCE TRANSLATIONAL MEDICINE | RESEARCH ARTICLE

HUNTINGTON'S DISEASE

## Evaluation of mutant huntingtin and neurofilament proteins as potential markers in Huntington's disease

Lauren M. Byrne[1]*[†], Filipe B. Rodrigues[1][†], Eileanor B. Johnson[1], Peter A. Wijeratne[2], Enrico De Vita[3,4], Daniel C. Alexander[2,5], Giuseppe Palermo[6], Christian Czech[6], Scott Schobel[6], Rachael I. Scahill[1], Amanda Heslegrave[7], Henrik Zetterberg[7,8,9,10], Edward J. Wild[1]*

• Biofluid markers change before imaging and clinical markers

Simplest way is to take the stage that maximises the likelihood for each patient

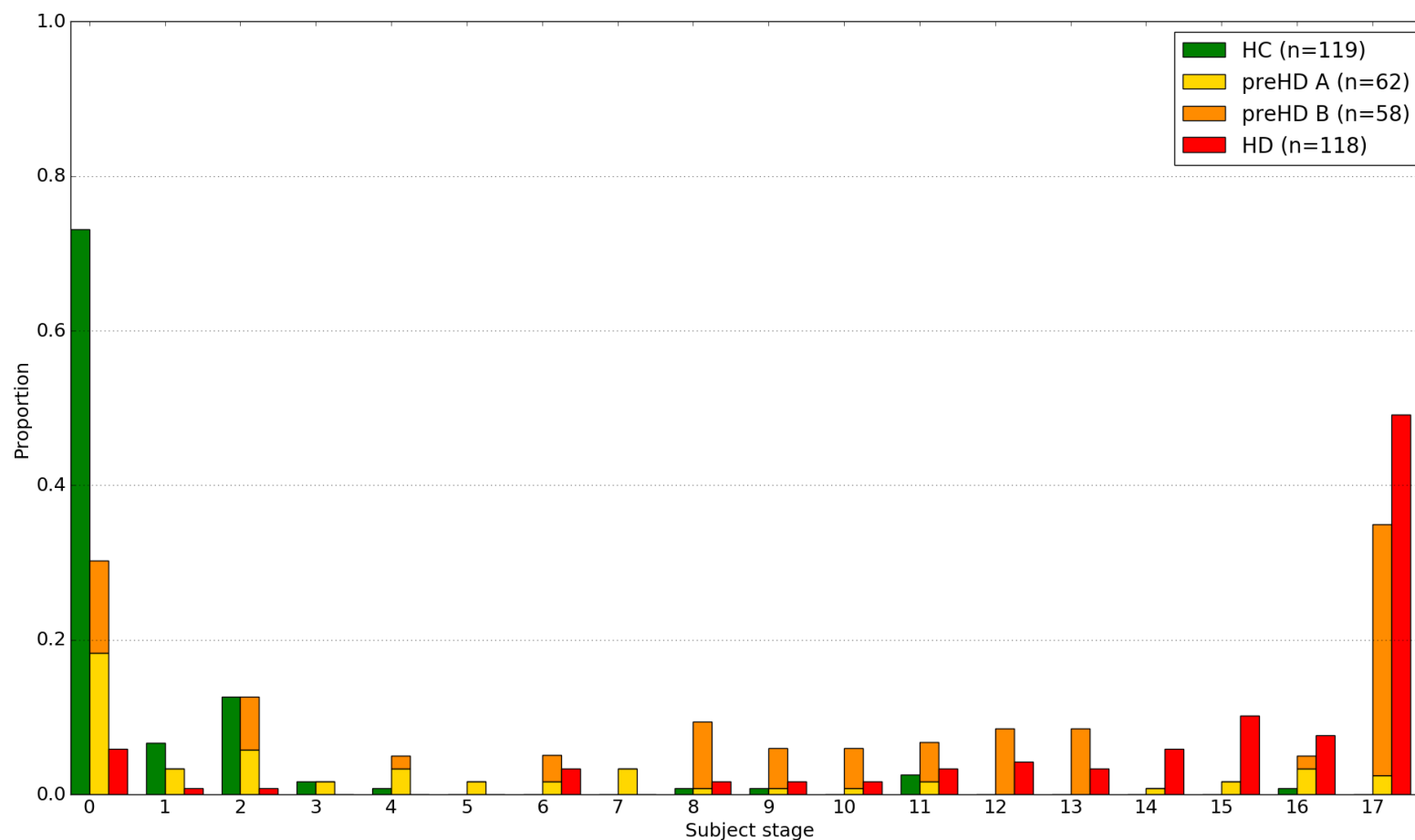$$argmax_k P(X_j | \bar{S}, k) = argmax_k P(k) \prod_{i=1}^{k} P(x_{ij}|E_i) \prod_{i=k+1}^{l} P(x_{ij}|\neg E_i)$$
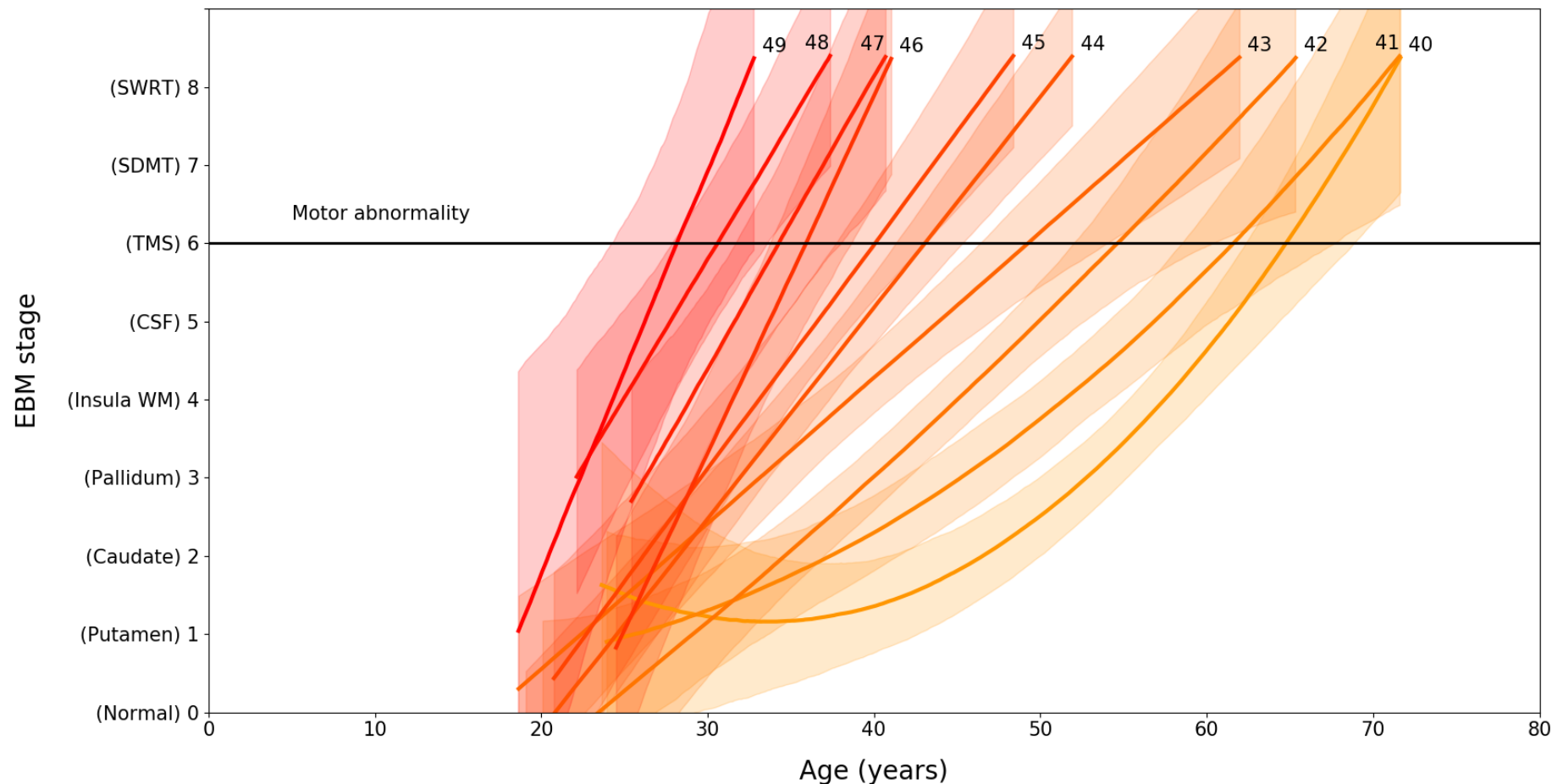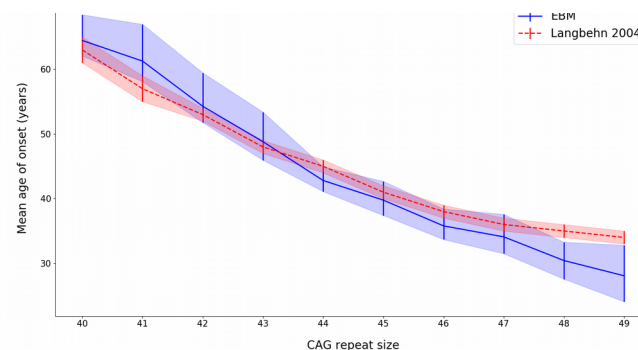
Simplest way is to take the stage that maximises the likelihood for each patient

$$argmax_k P(X_j | \bar{S}, k) = argmax_k P(k) \prod_{i=1}^{k} P(x_{ij} | E_i) \prod_{i=k+1}^{l} P(x_{ij} | \neg E_i)$$
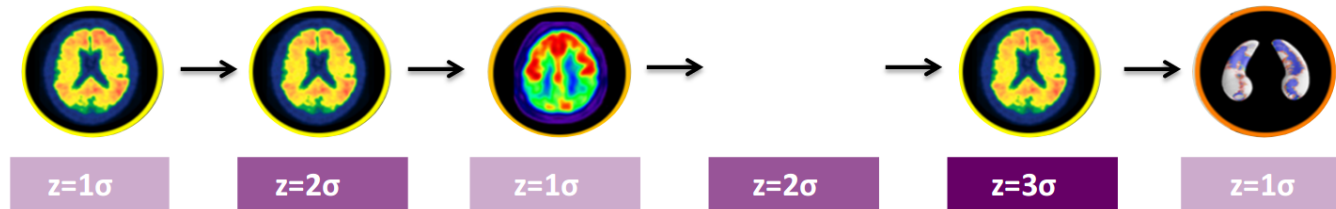


36

- Estimate age at event e.g.
  for CAG 40, WM atrophy at ~60 years old
  for CAG 49, WM atrophy at ~25 years old

- Age of onset agrees well with gold standard

37

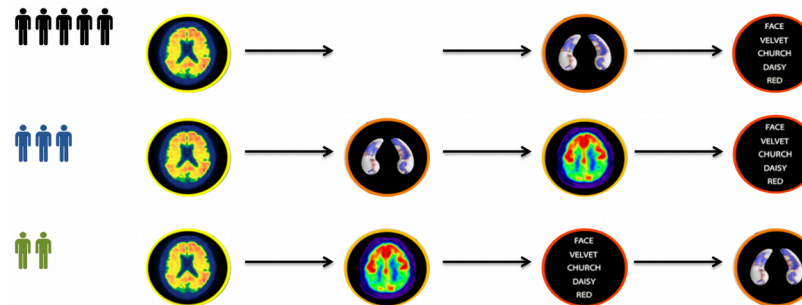1. Continuous generalisation of EBM: instead of instantaneous abnormality, markers are a linear combination of z-scores
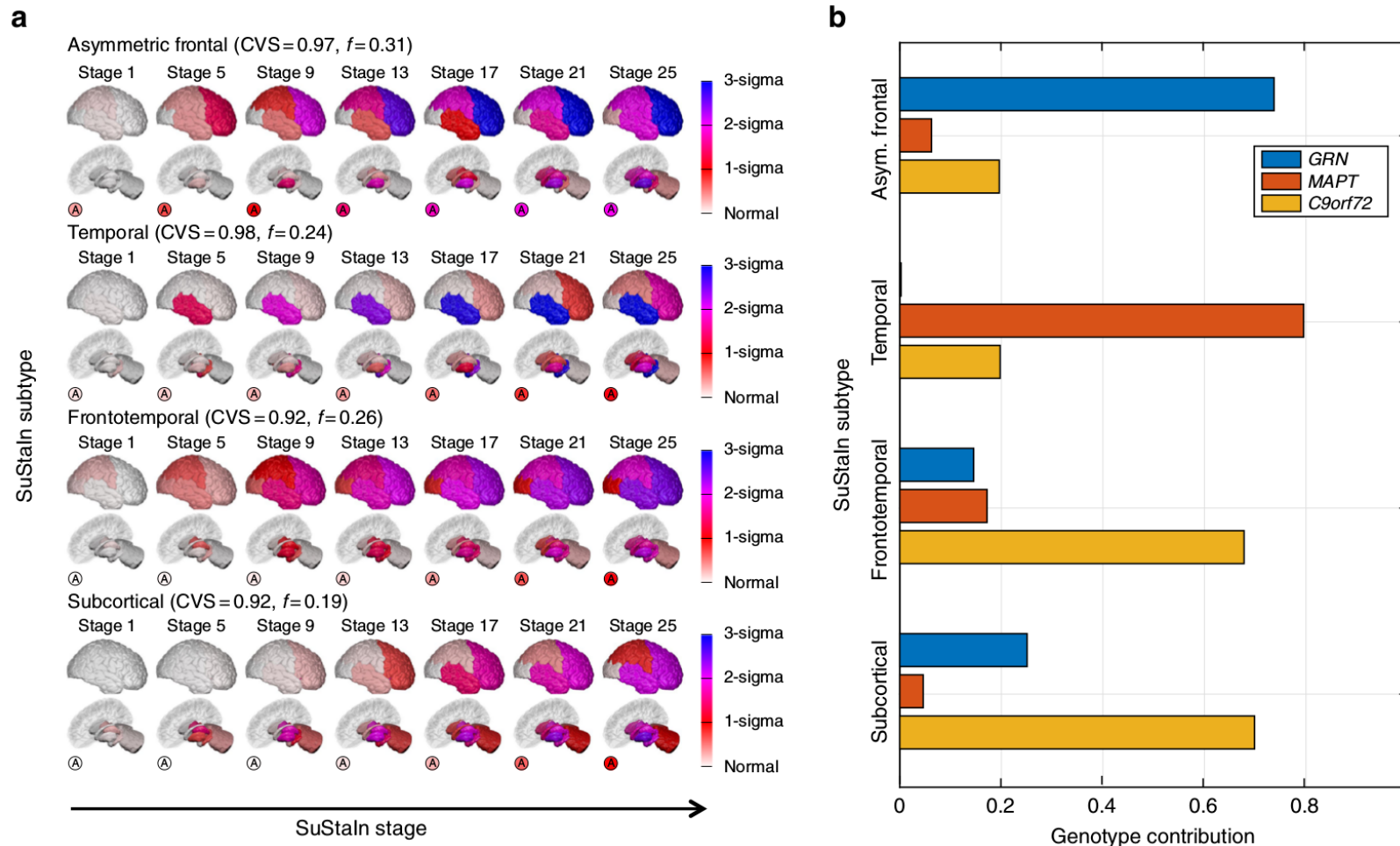


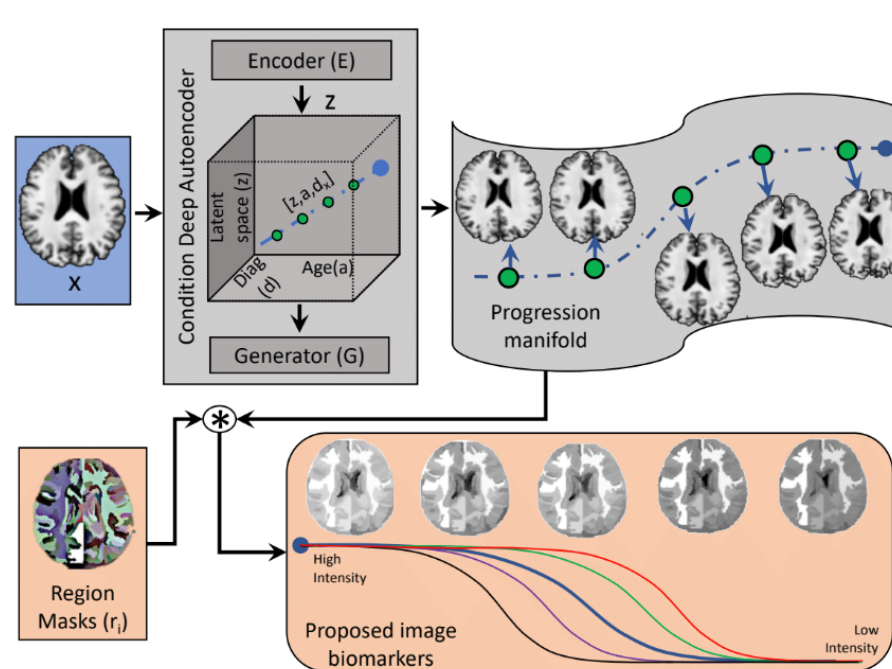| z=1σ | z=2σ | z=1σ | z=2σ | z=3σ | z=1σ |

"Z-score model"

2. Total model is mixture of linear z-score models: grouped into clusters with distinct progression patterns
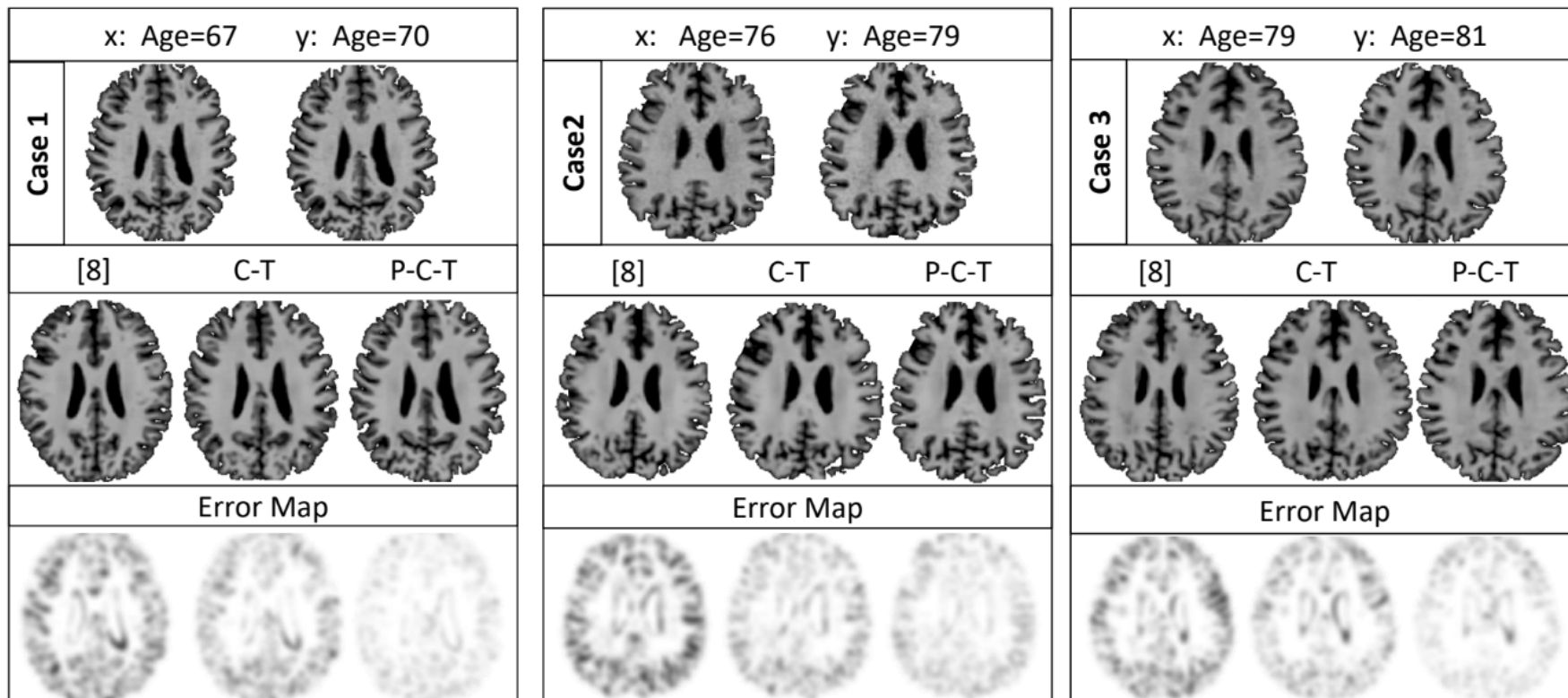


"Algorithm"

Gain this extra information just by generalising event-based model
– pretty neat

Deep learning disease trajectories using generative adversarial networks
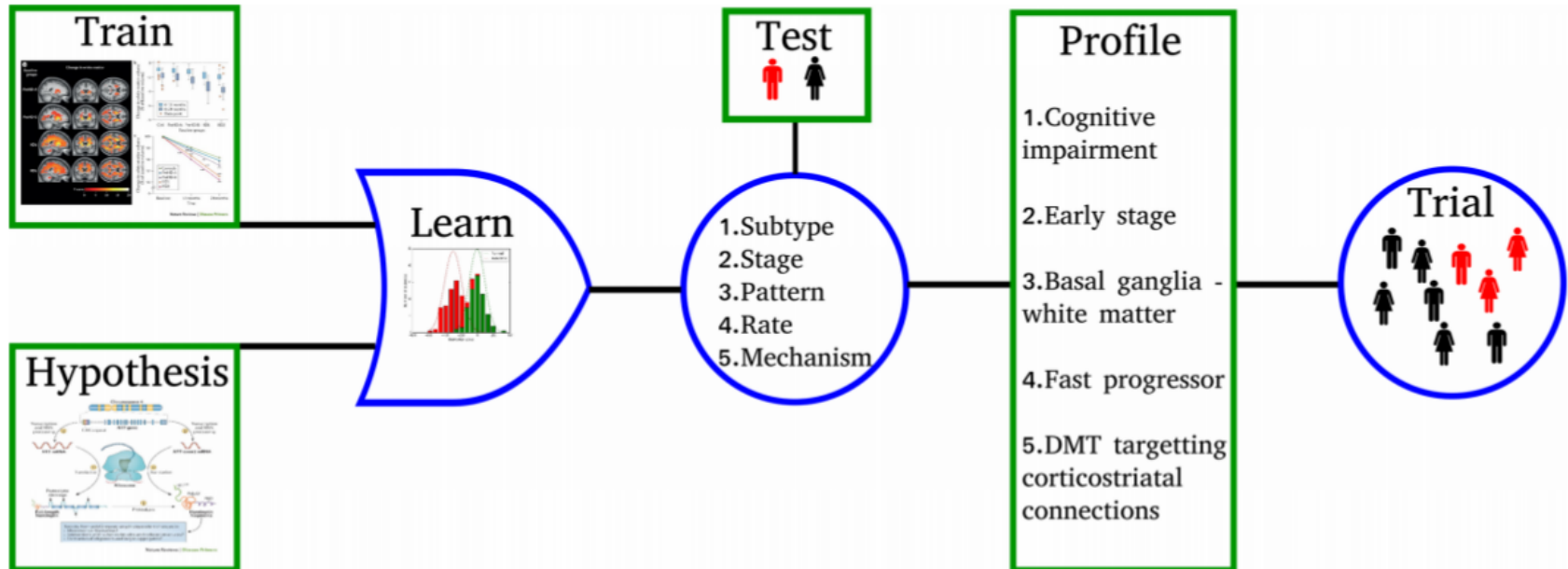
- also used in HEP e.g. CaloGAN, Paganini, Oliveira, Nachman. 2017.

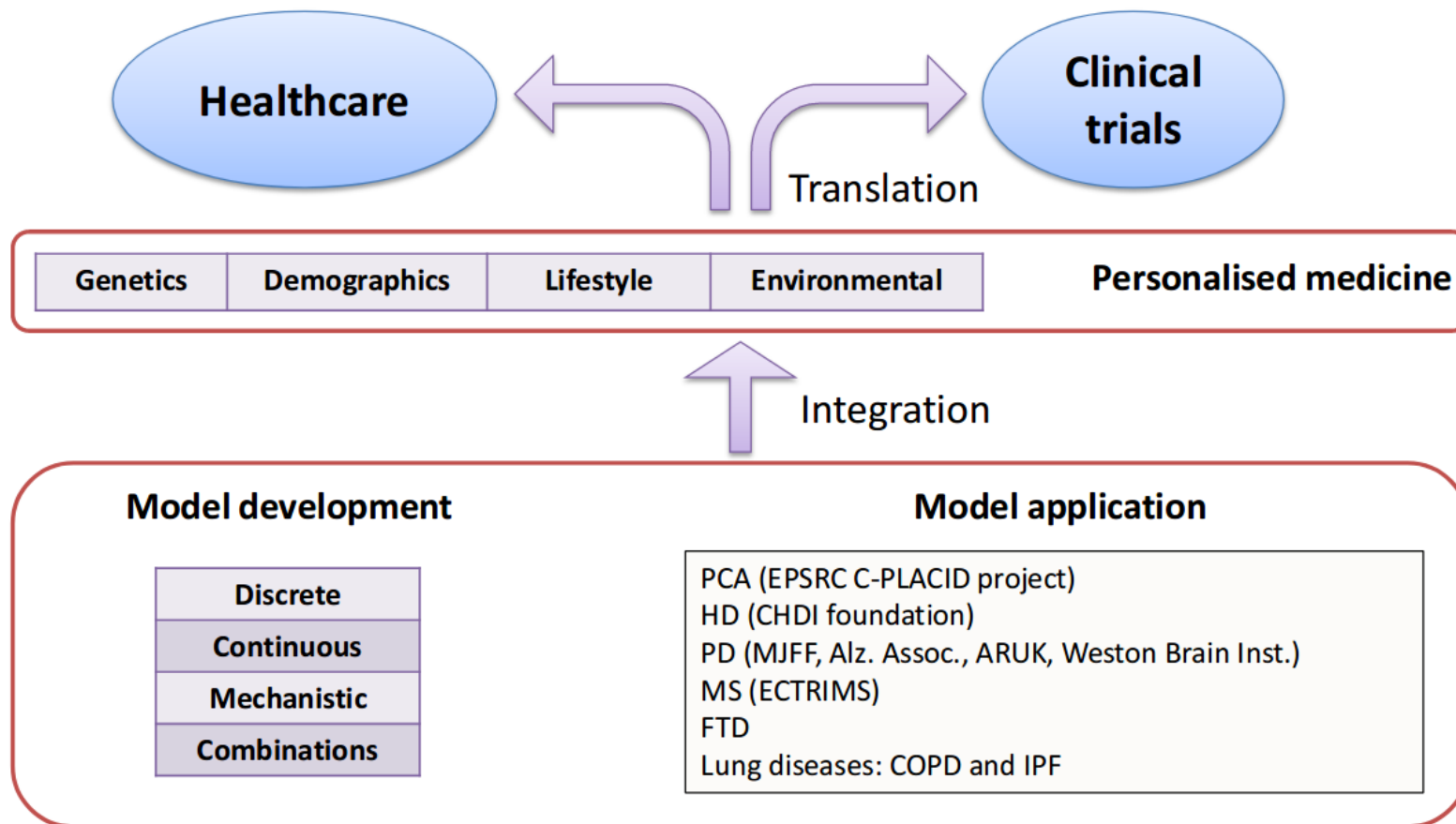Deep learning disease trajectories using generative adversarial networks

- also used in HEP e.g. CaloGAN, Paganini, Oliveira, Nachman. 2017.

Patient data + machine learning = personalised profiles for clinical trial design



Model can be used for both prospective and retrospective analysis

→ Save money and time
→ Optimise trial design

• Presented computational methods to extract information from large and varied datasets

• Machine learning methods are suitable for medical problems – i.e. inferring patterns from complex systems

• Still much to do – can we understand the mechanisms themselves?

• What can HEP and CS learn from each other?



https://www.slideshare.net/mlreview/tutorial-on-deep-generative-models