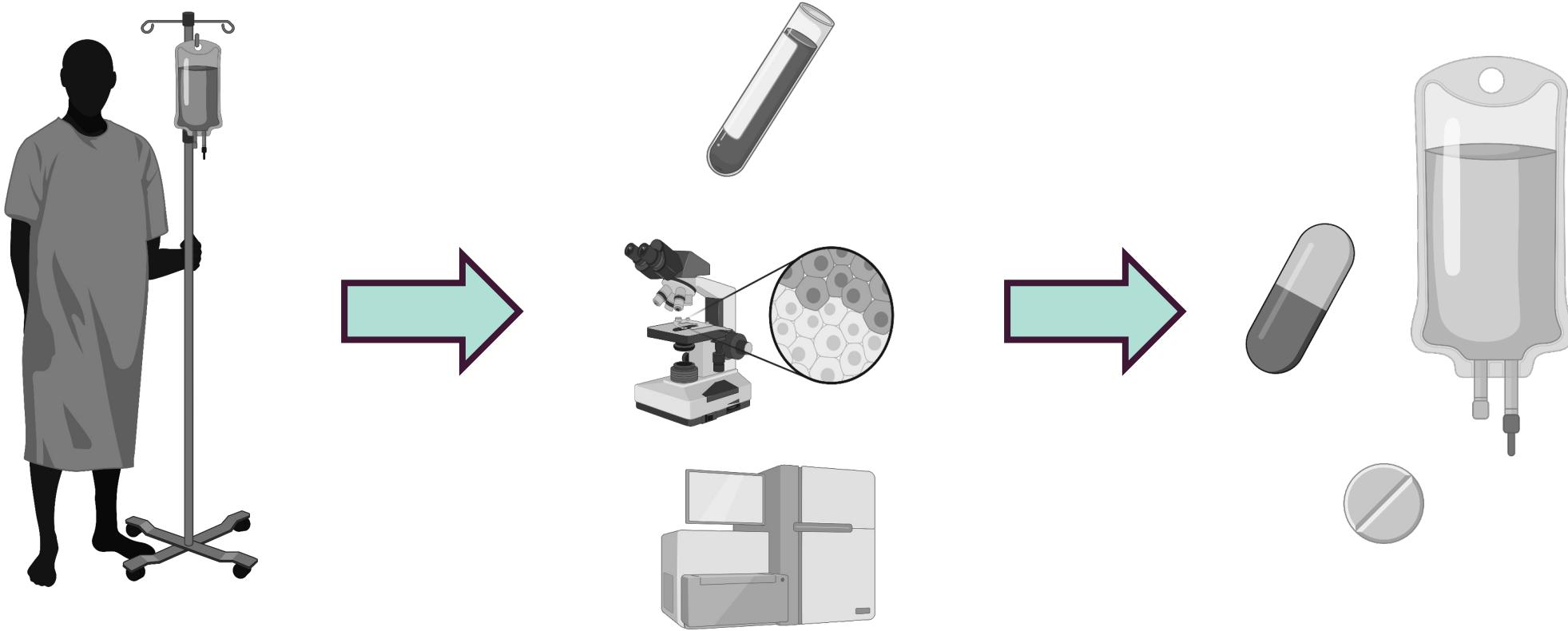


Live brain cancer classification during surgery

Basic Machine Learning Bioinformatics
09-03-2023

Personalized medicine

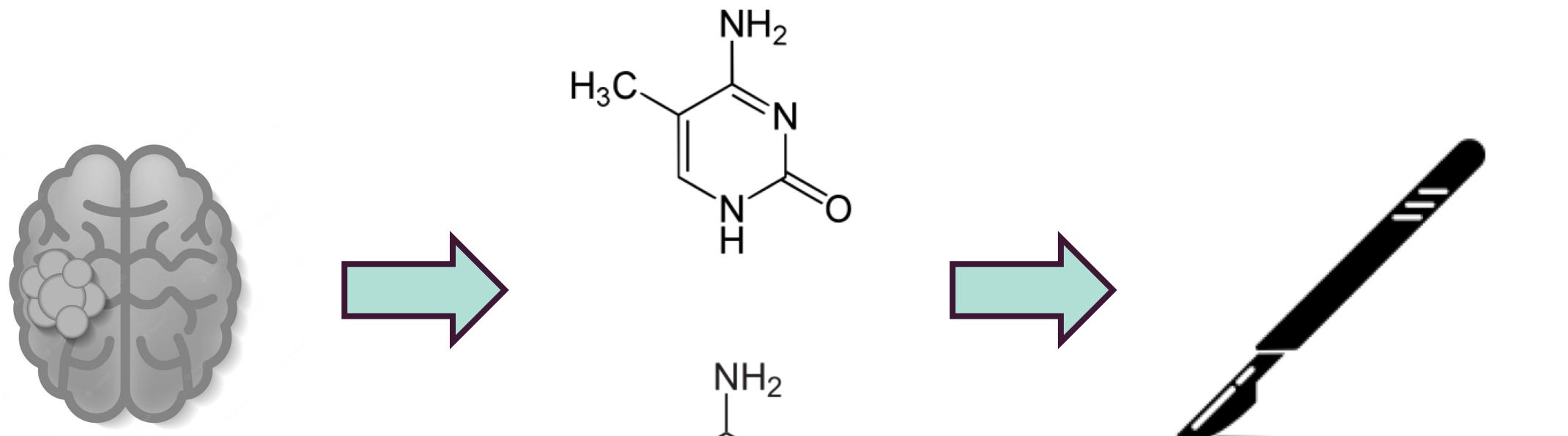


Patient

Biomarker

Treatment

Personalized medicine



**Brain
cancer**

**Methylation
status**

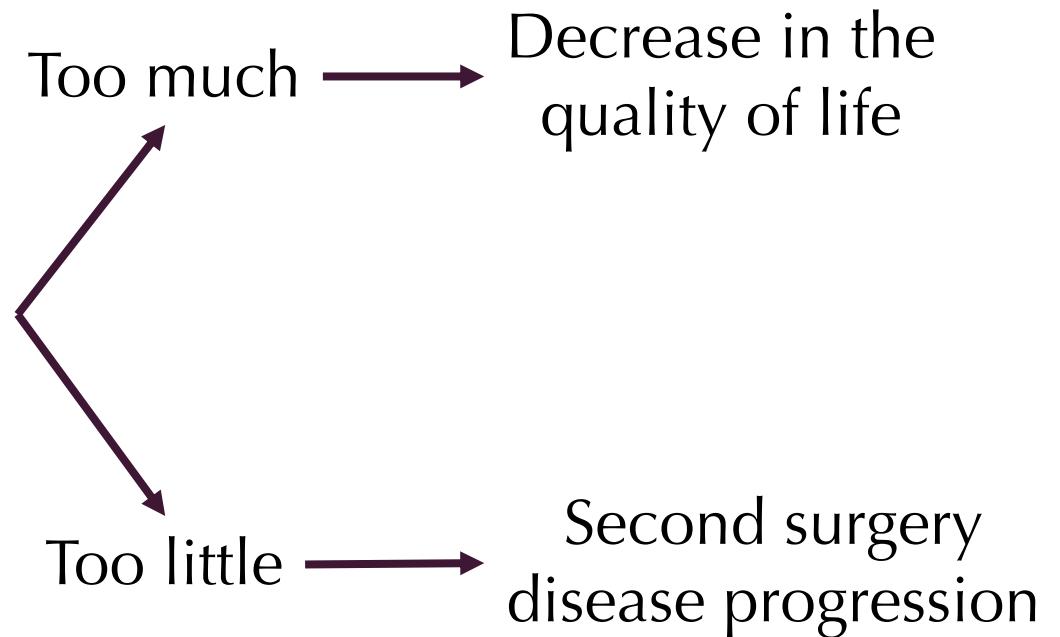
**Personalized
surgery**

Personalized surgery

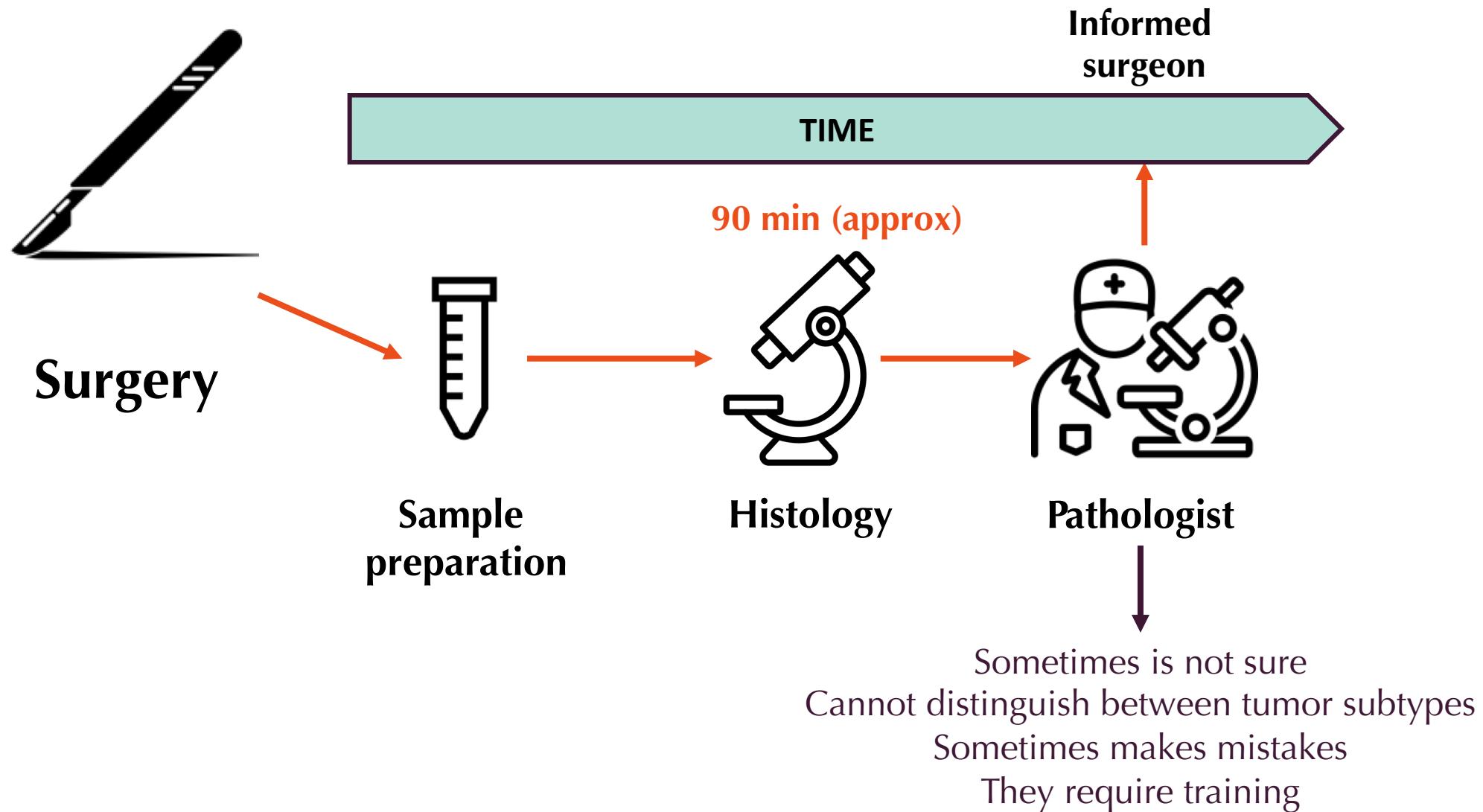


Personalized
surgery

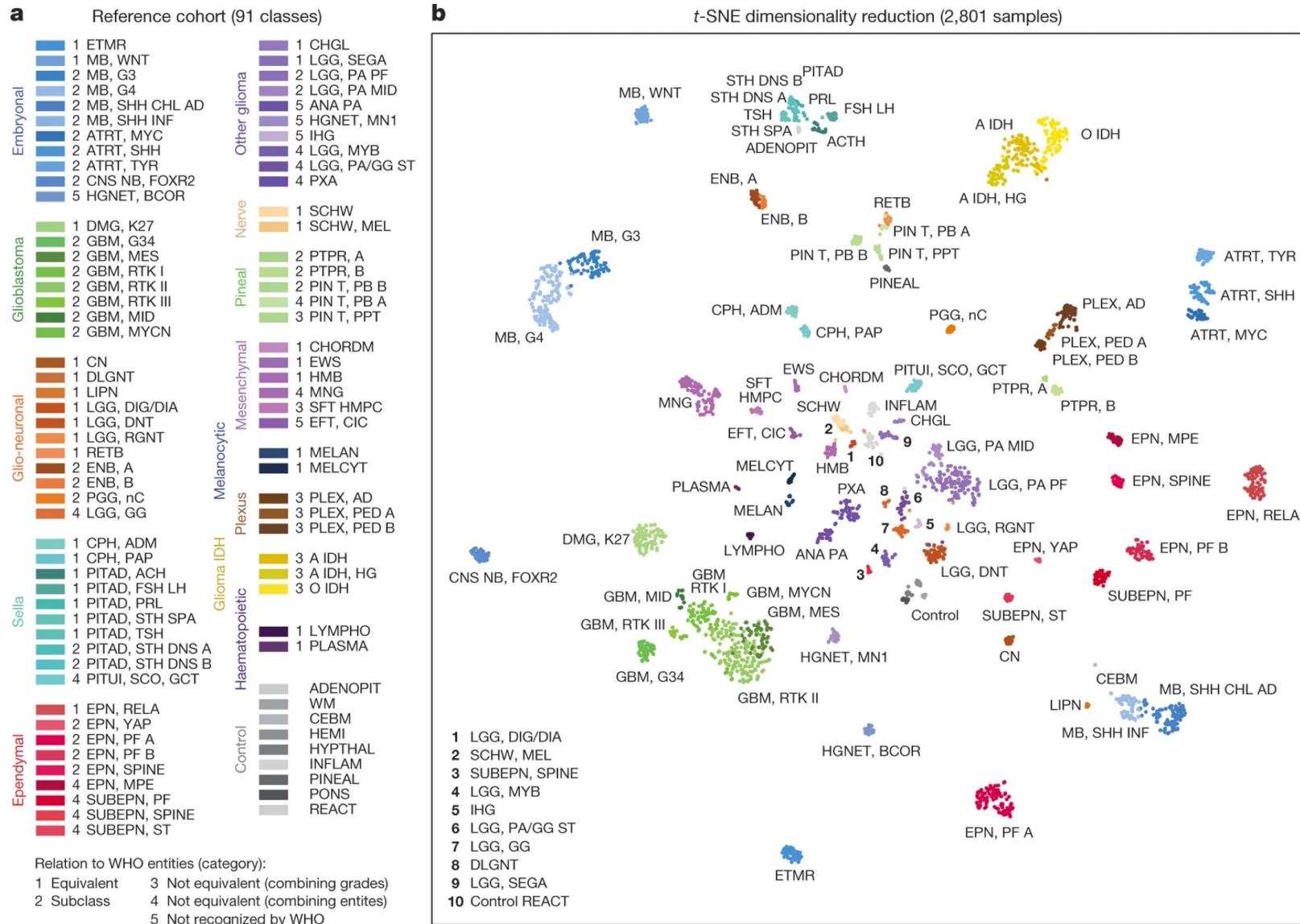
How much resection
should be done?



Peri-operative CNS classification



Methylation as CNS class biomarker

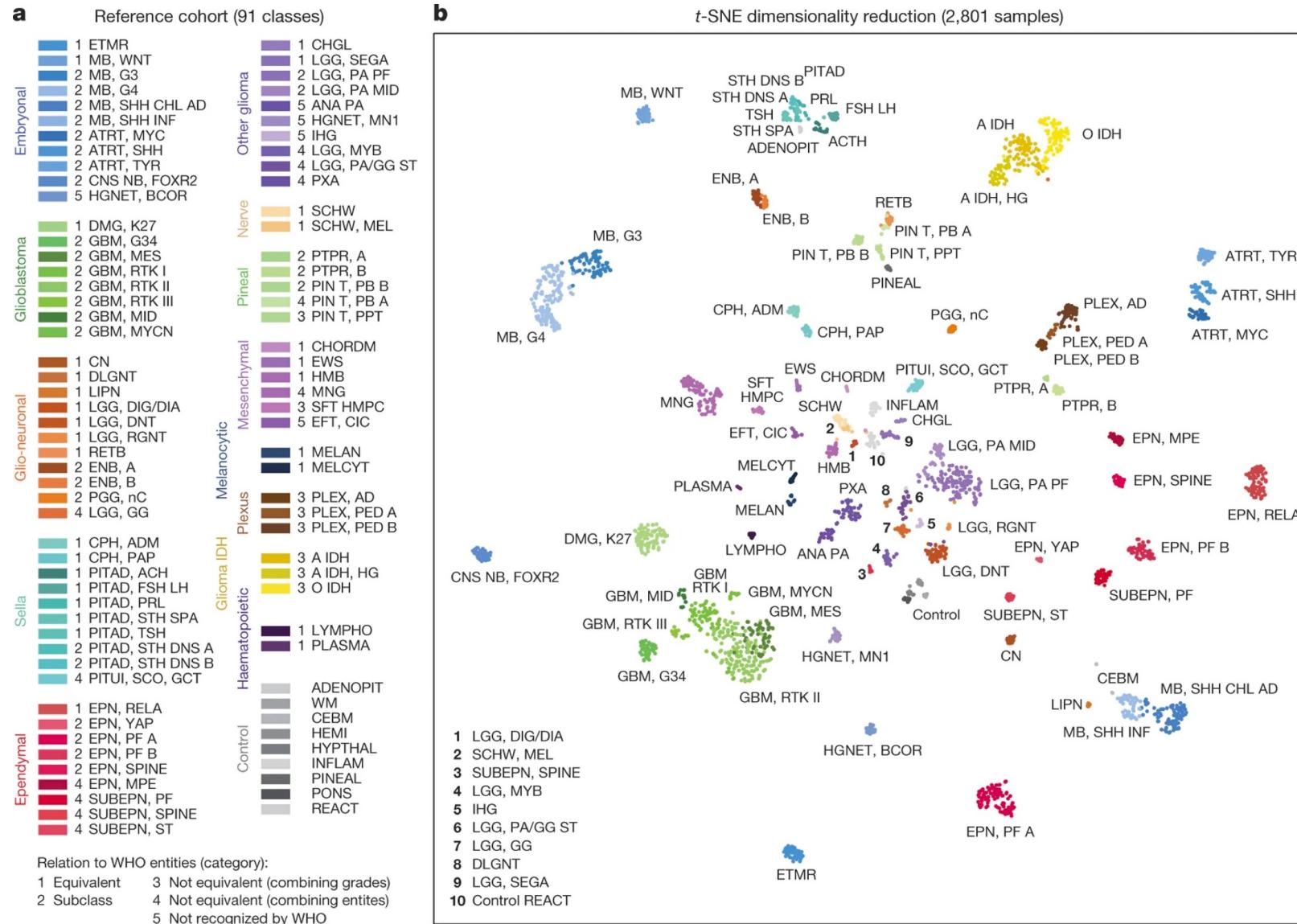


2801 CNS samples
82 tumor classes
9 control classes

Illumina Infinium array

WHAT DO YOU SEE?

Methylation as CNS class biomarker



2801 CNS samples
82 tumor classes
9 control classes

Illumina Infinium array

Methylation is a really good biomarker
Unsupervised clustering already gives very clear clusters

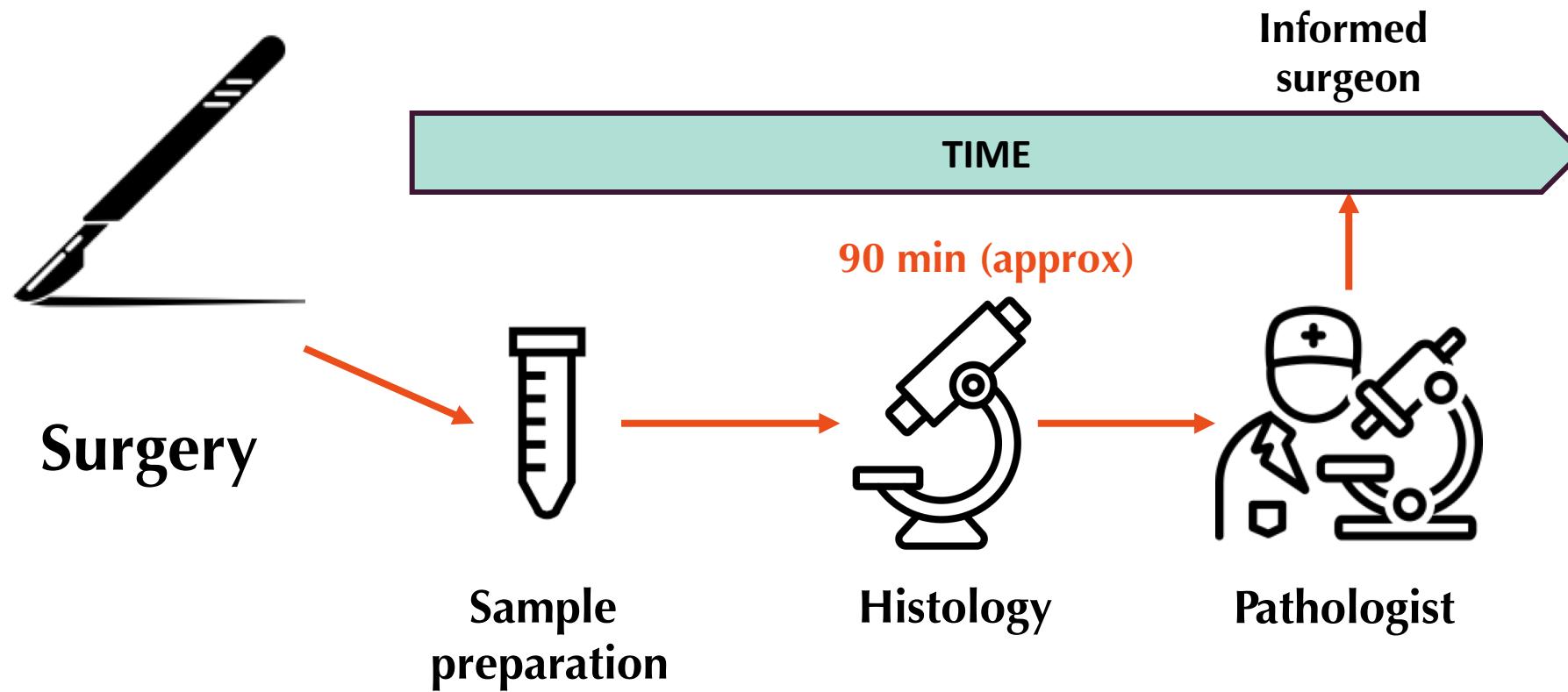
Illumina methylation arrays



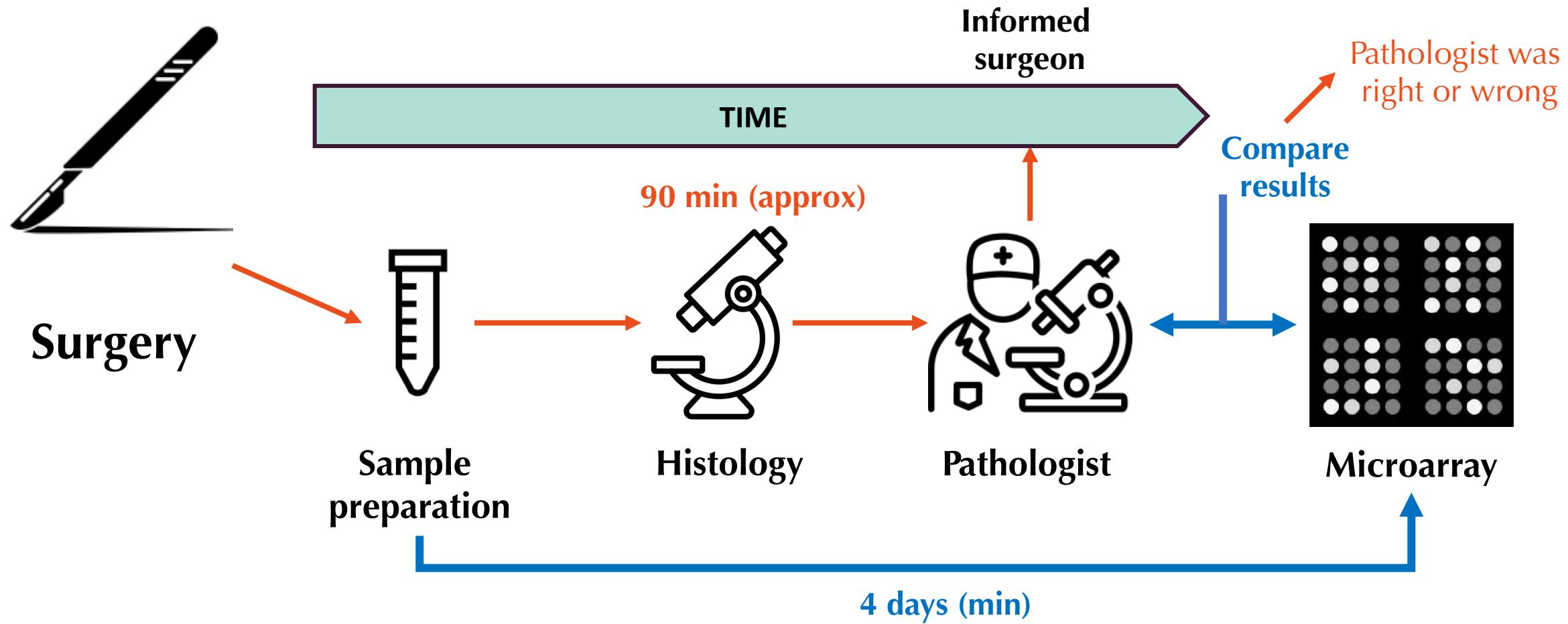
Illumina Infinium array

- **450k CpG methylation sites**
- **8 samples simultaneously**
- **Genome-wide coverage**
- **4 days workflow**

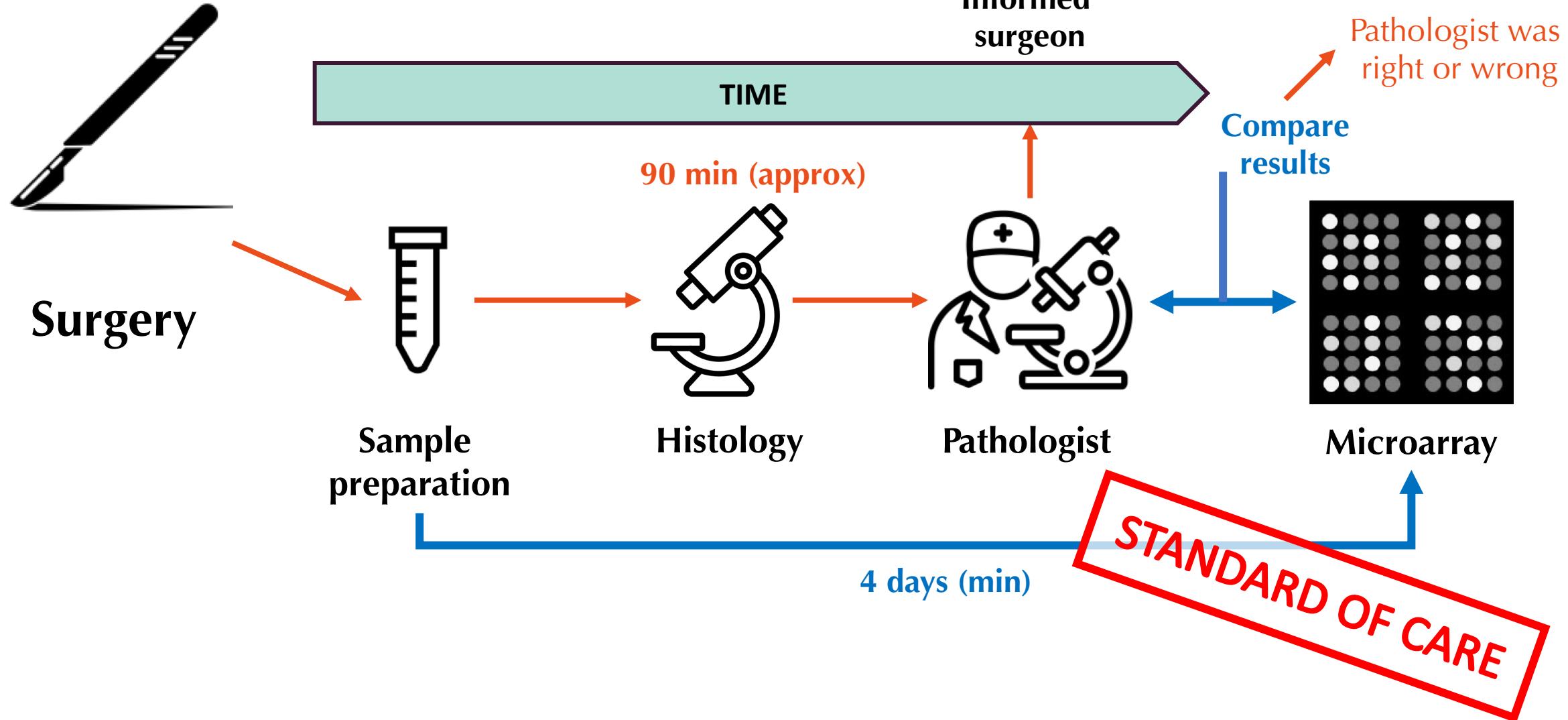
Peri-operative CNS classification



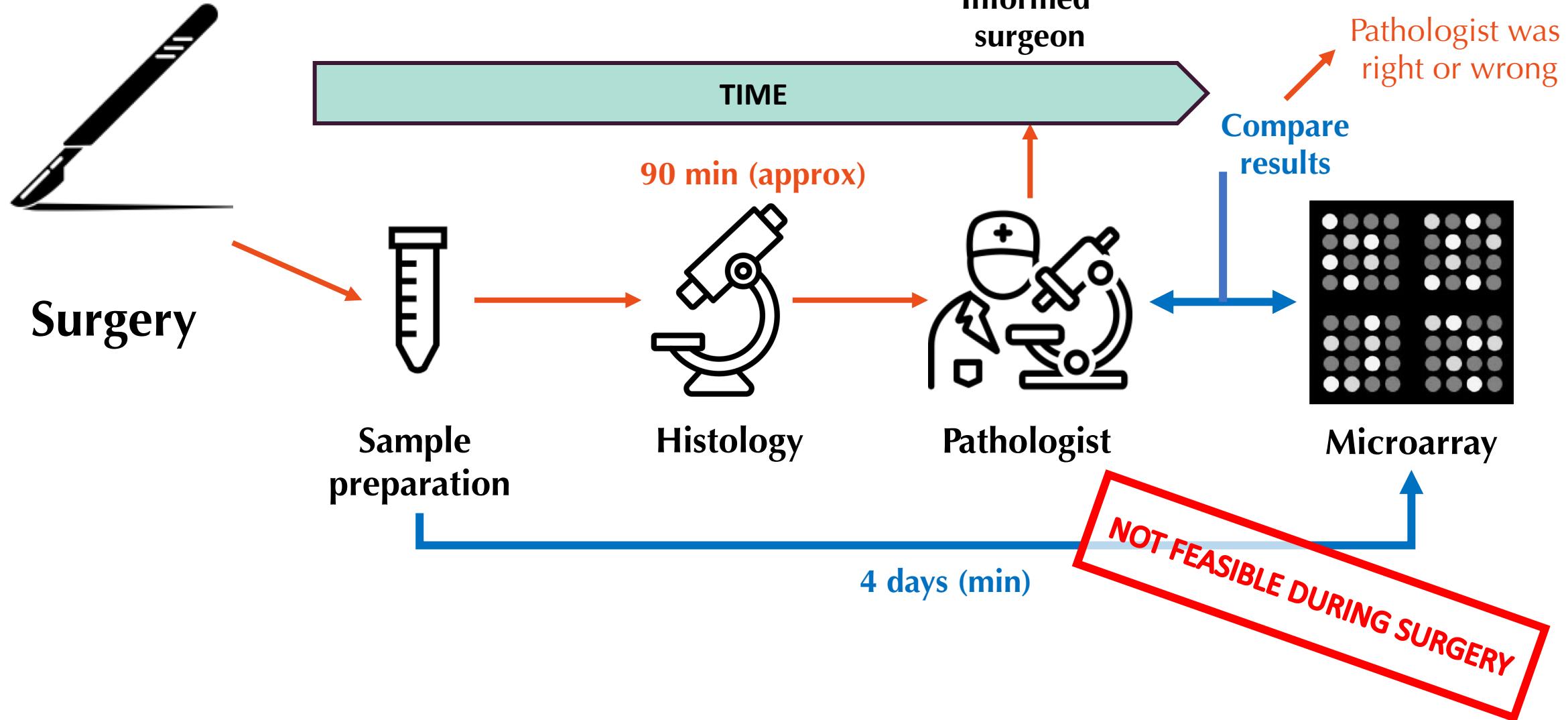
Peri/post-operative CNS classification



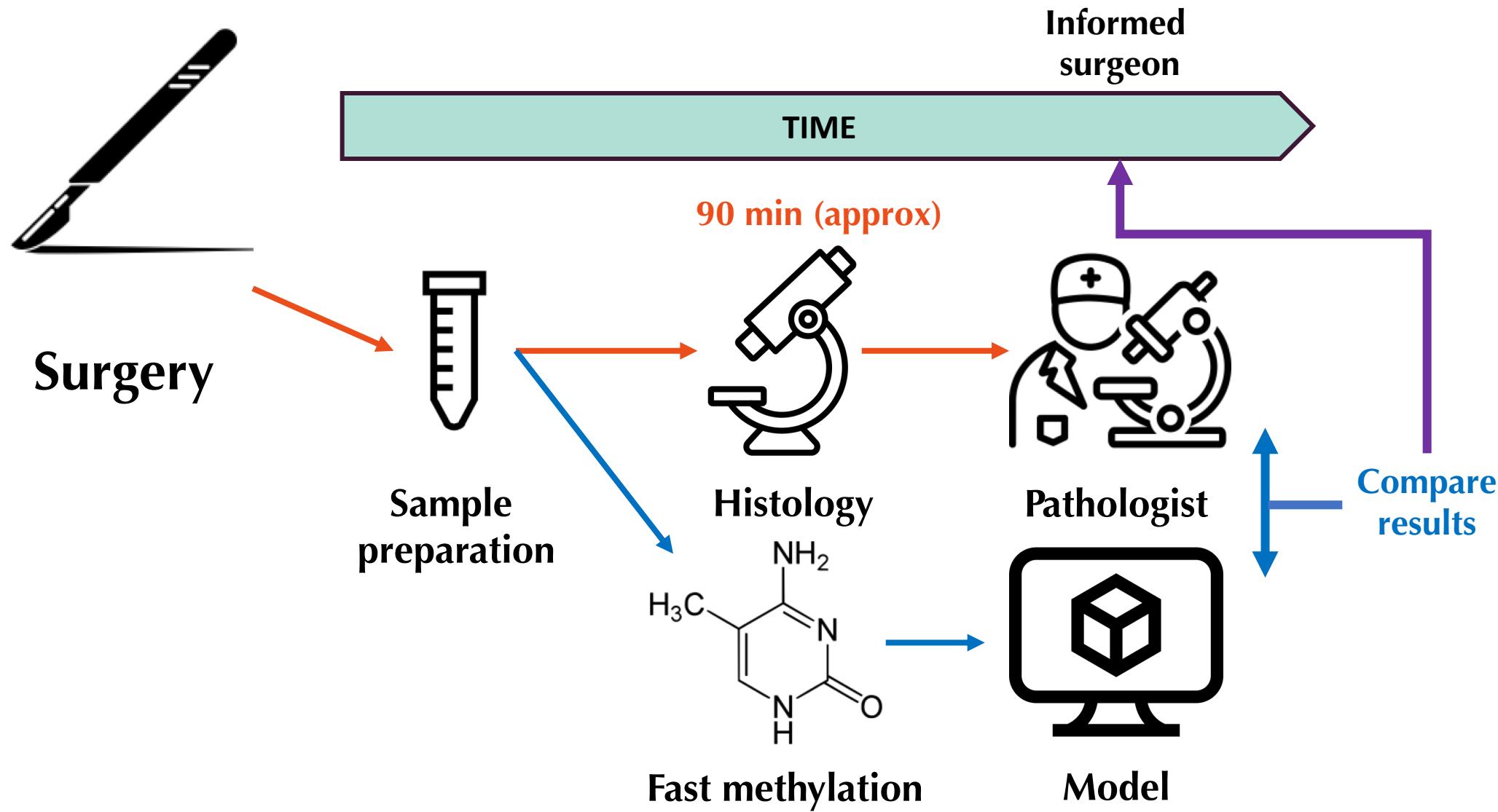
Peri/post-operative CNS classification



Peri/post-operative CNS classification



Peri-operative CNS classification



Fast methylation detection



Nanopore sequencing

- Native DNA sequencing
- Potential genome-wide coverage
- 15 minutes sample preparation
- Real time data output

Sparsity problem

Microarray

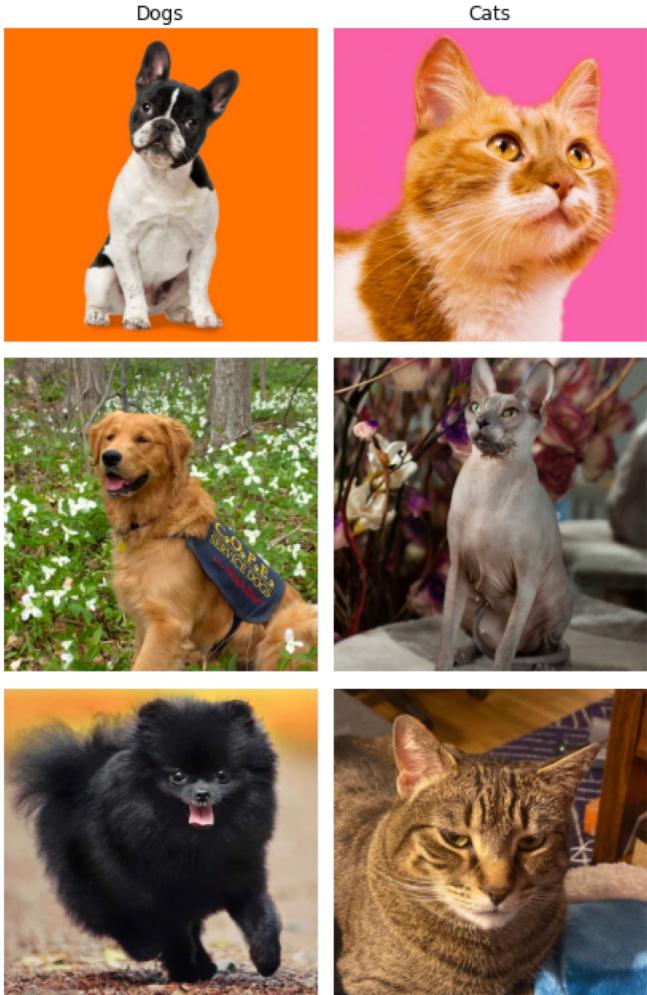
0% missing values
Continuous signal



Sparsity problem

Microarray

0% missing values
Continuous signal



Nanopore

>90% missing values
Binary signal



Unknown which missing values

Sparsity problem

TRAIN

Microarray

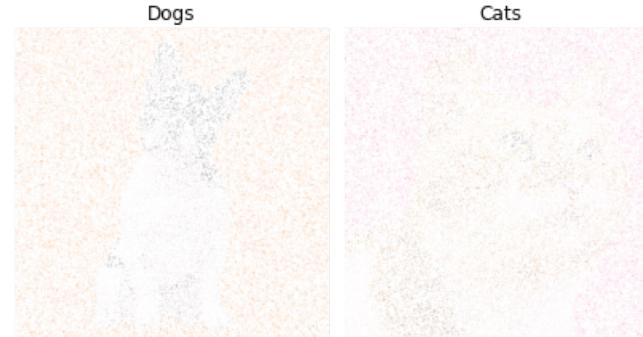
0% missing values
Continuous signal



TEST

Nanopore

>90% missing values
Binary signal



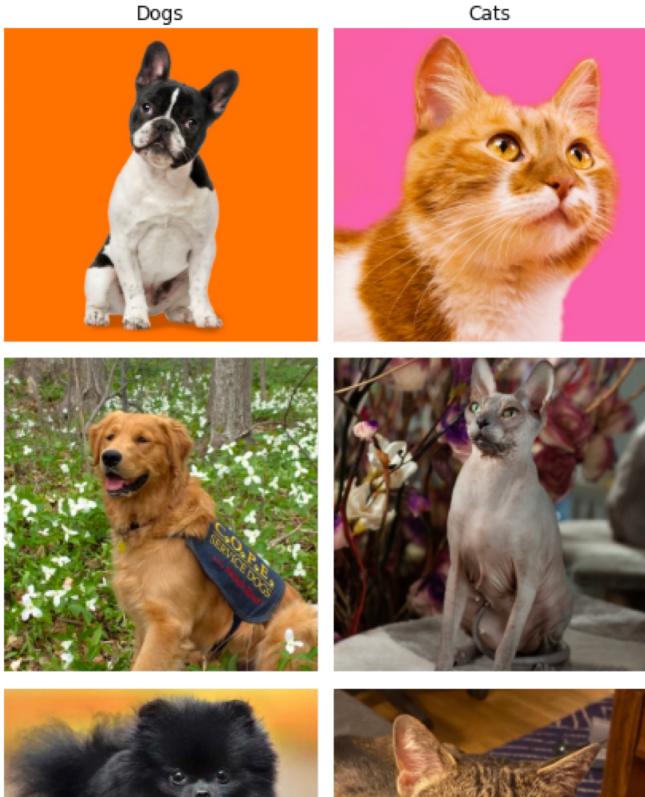
Unknown which missing values

Sparsity problem

TRAIN

Microarray

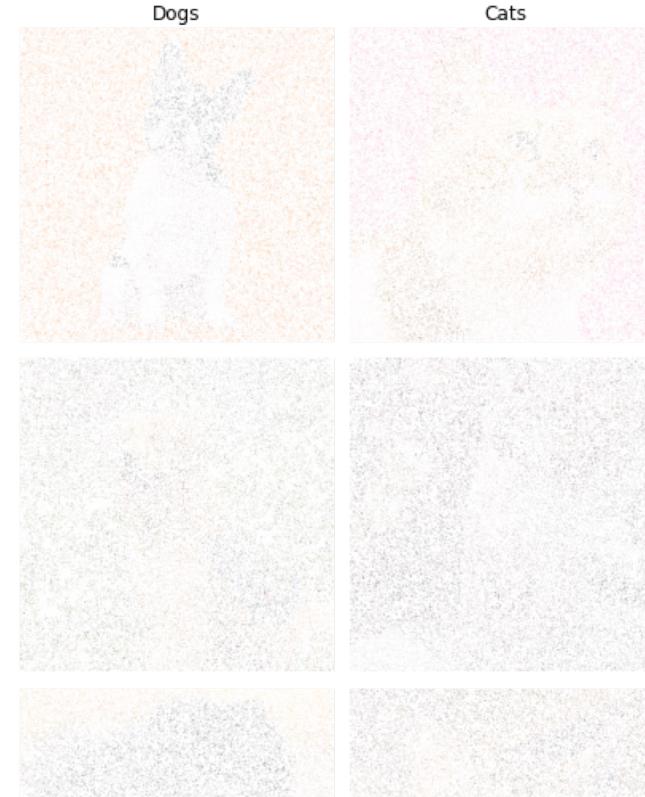
0% missing values
Continuous signal



TEST

Nanopore

>90% missing values
Binary signal



Can we develop a method that can handle an arbitrary random amount of missing values?

Sparsity problem

TRAIN

Microarray

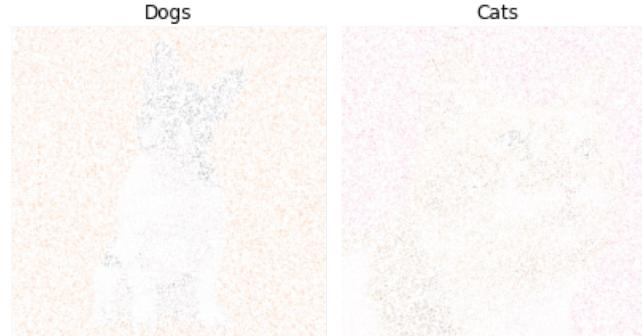
0% missing values
Continuous signal



TEST

Nanopore

>90% missing values
Binary signal



ANY IDEAS?

Existing approaches

ORIGINAL ARTICLE

Neuropathology and
Applied Neurobiology
JOURNAL OF THE BRITISH NEUROPATHOLOGICAL SOCIETY

WILEY

Robust methylation-based classification of brain tumours using nanopore sequencing

Luis P. Kuschel¹  | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Masliah-Planchon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaih⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9} 

Robust methylation-based classification of brain tumours using nanopore sequencing

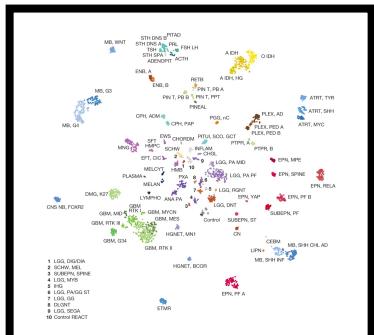
Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbally⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}



Nanopore sequencing



2801x450000



Microarray data

Robust methylation-based classification of brain tumours using nanopore sequencing

Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaili⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}

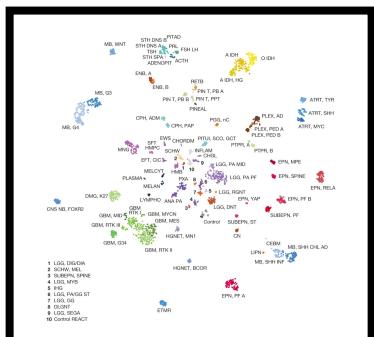
Existing approaches



Nanopore sequencing



2801x450000



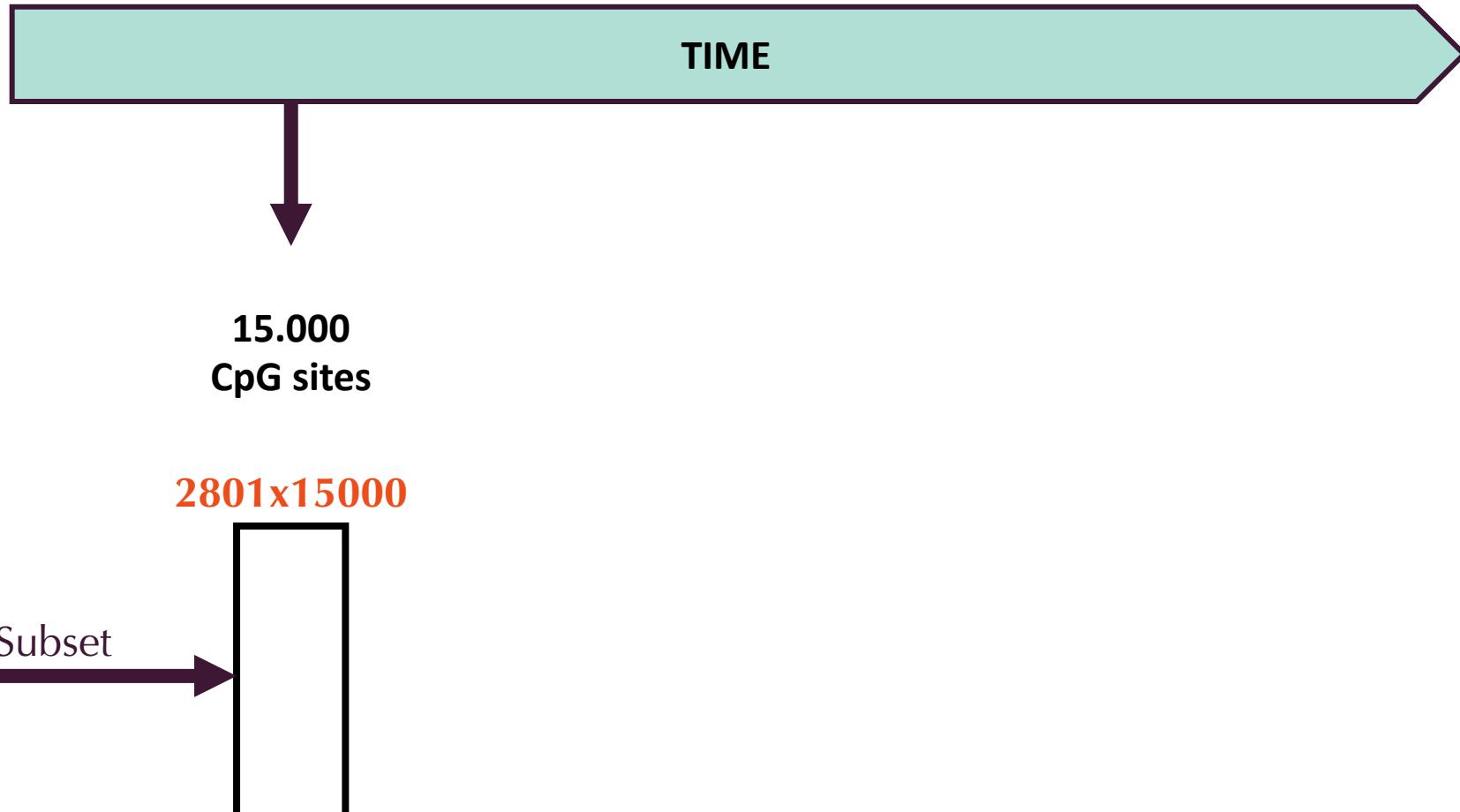
Microarray data

Robust methylation-based classification of brain tumours using nanopore sequencing

Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaili⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}



Nanopore sequencing



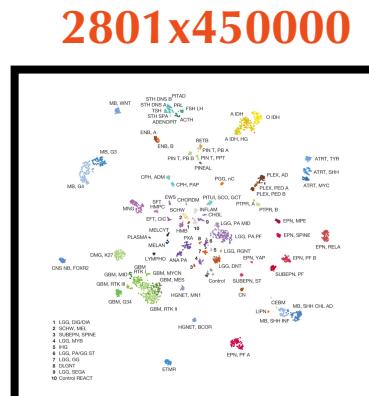
Microarray data

Robust methylation-based classification of brain tumours using nanopore sequencing

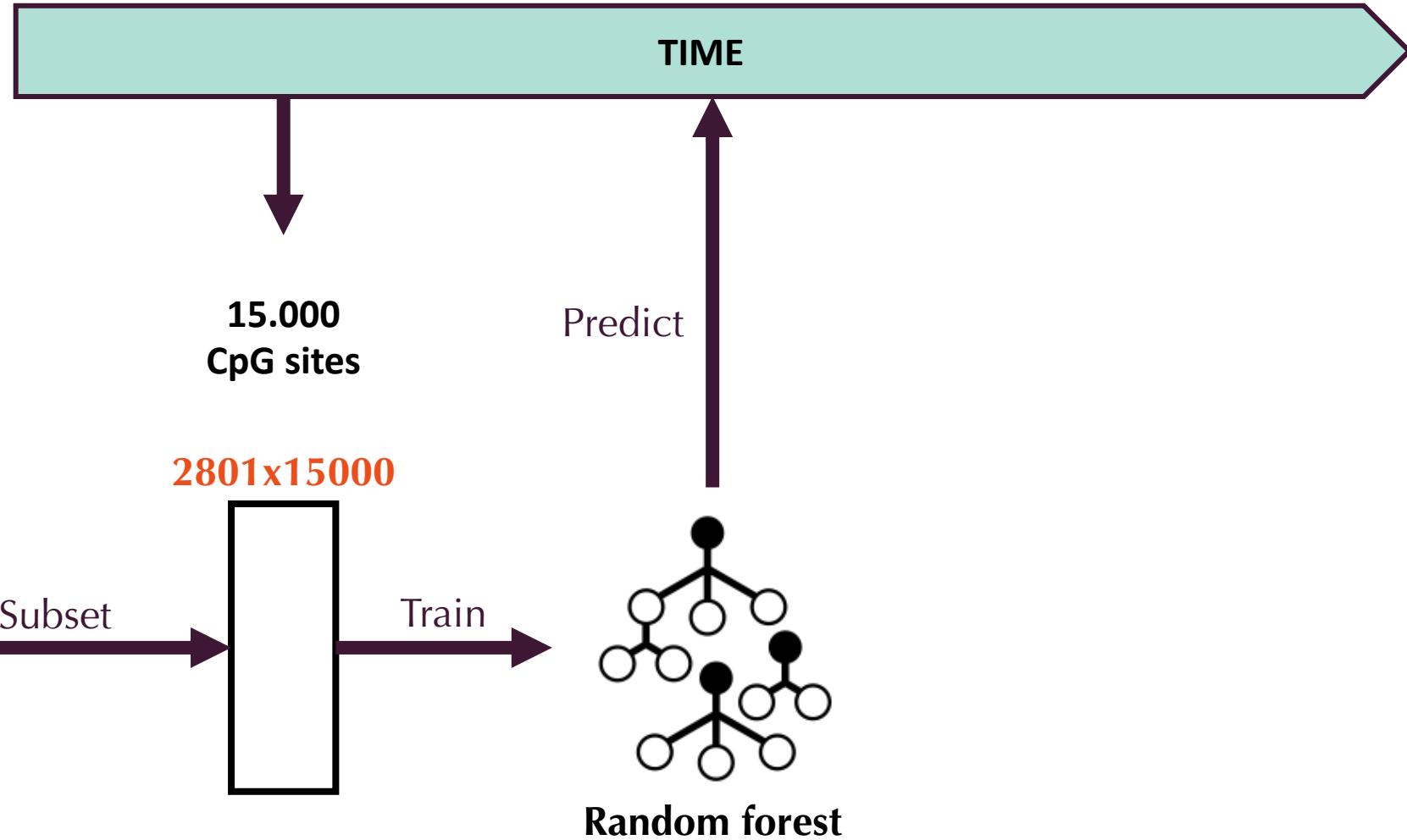
Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaili⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}



Nanopore sequencing



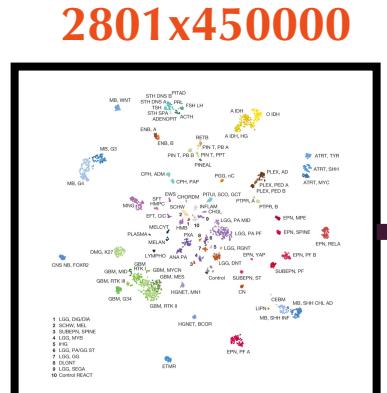
Microarray data



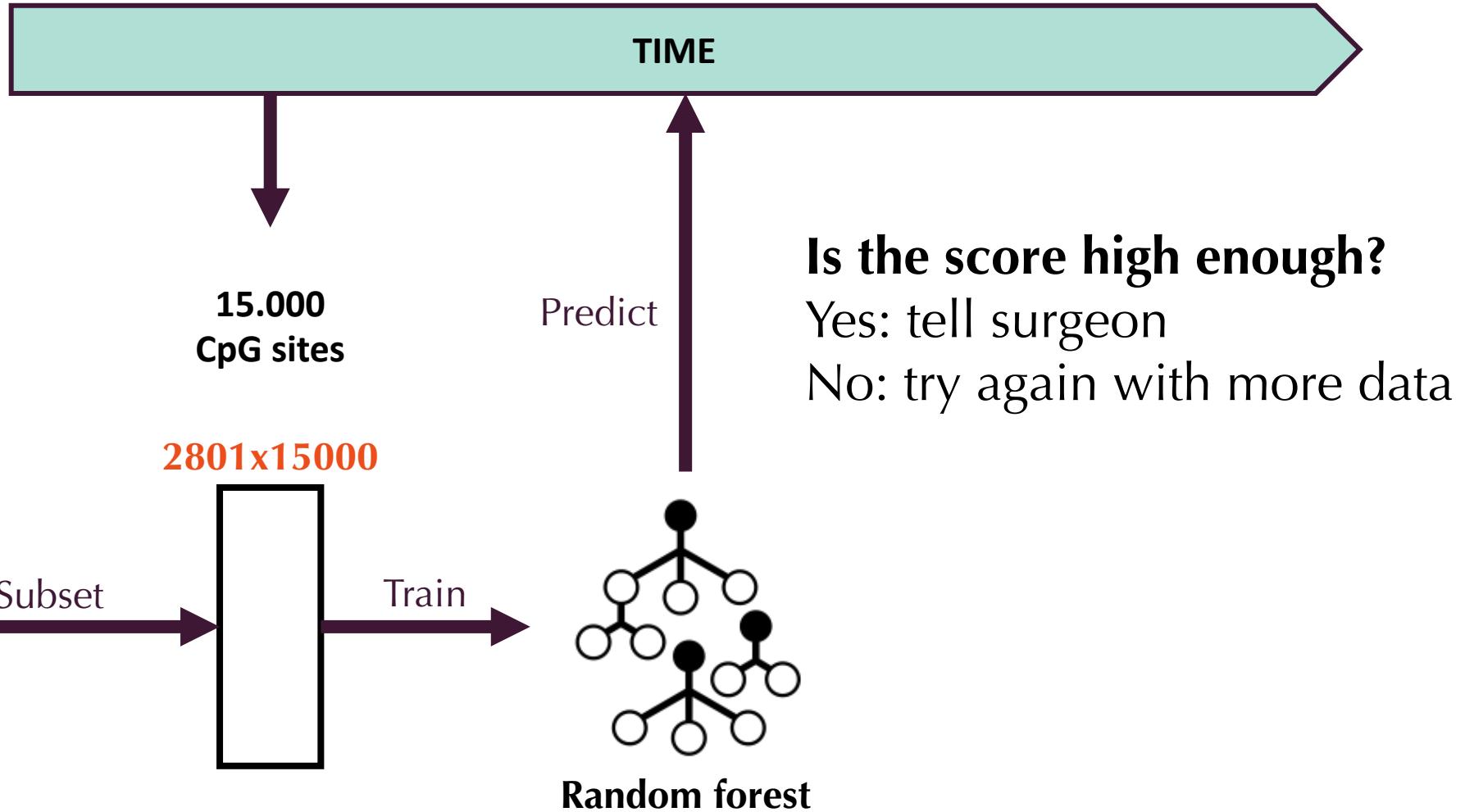
Existing approaches



Nanopore sequencing



Microarray data



Robust methylation-based classification of brain tumours using nanopore sequencing

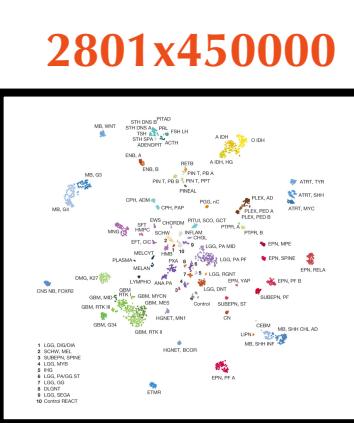
Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaili⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}

Robust methylation-based classification of brain tumours using nanopore sequencing

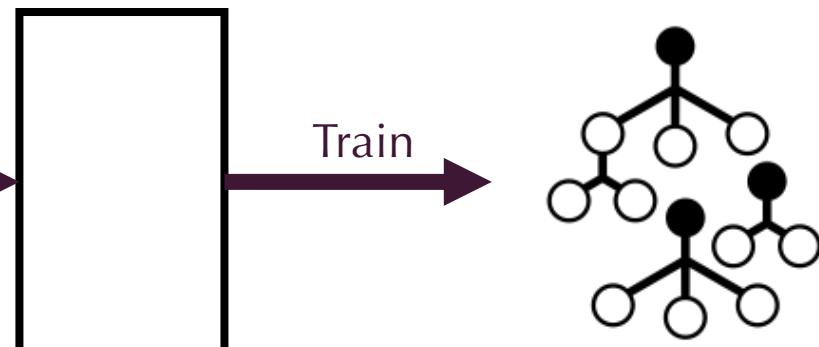
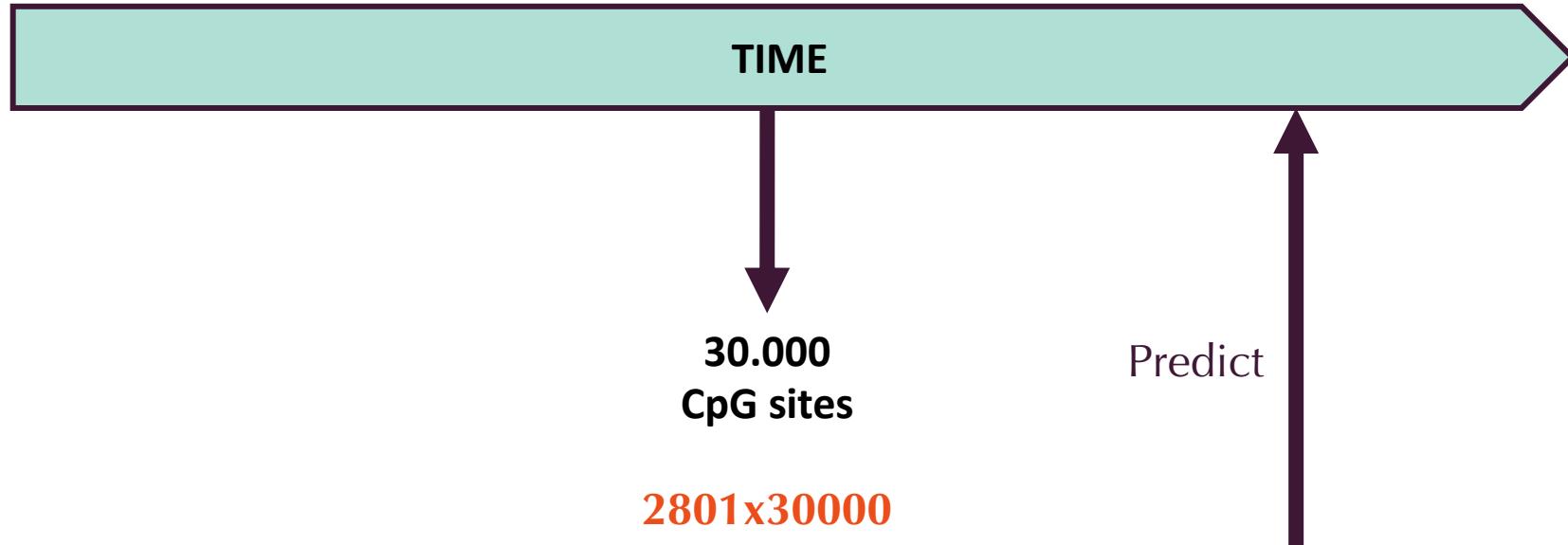
Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Masliah-Planchet³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbally⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}



Nanopore sequencing



Existing approaches



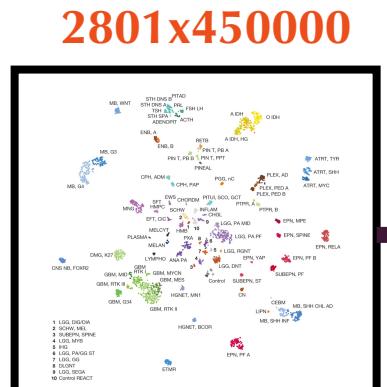
Robust methylation-based classification of brain tumours using nanopore sequencing

Luis P. Kuschel¹ | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Maslah-Planckon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaili⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9}

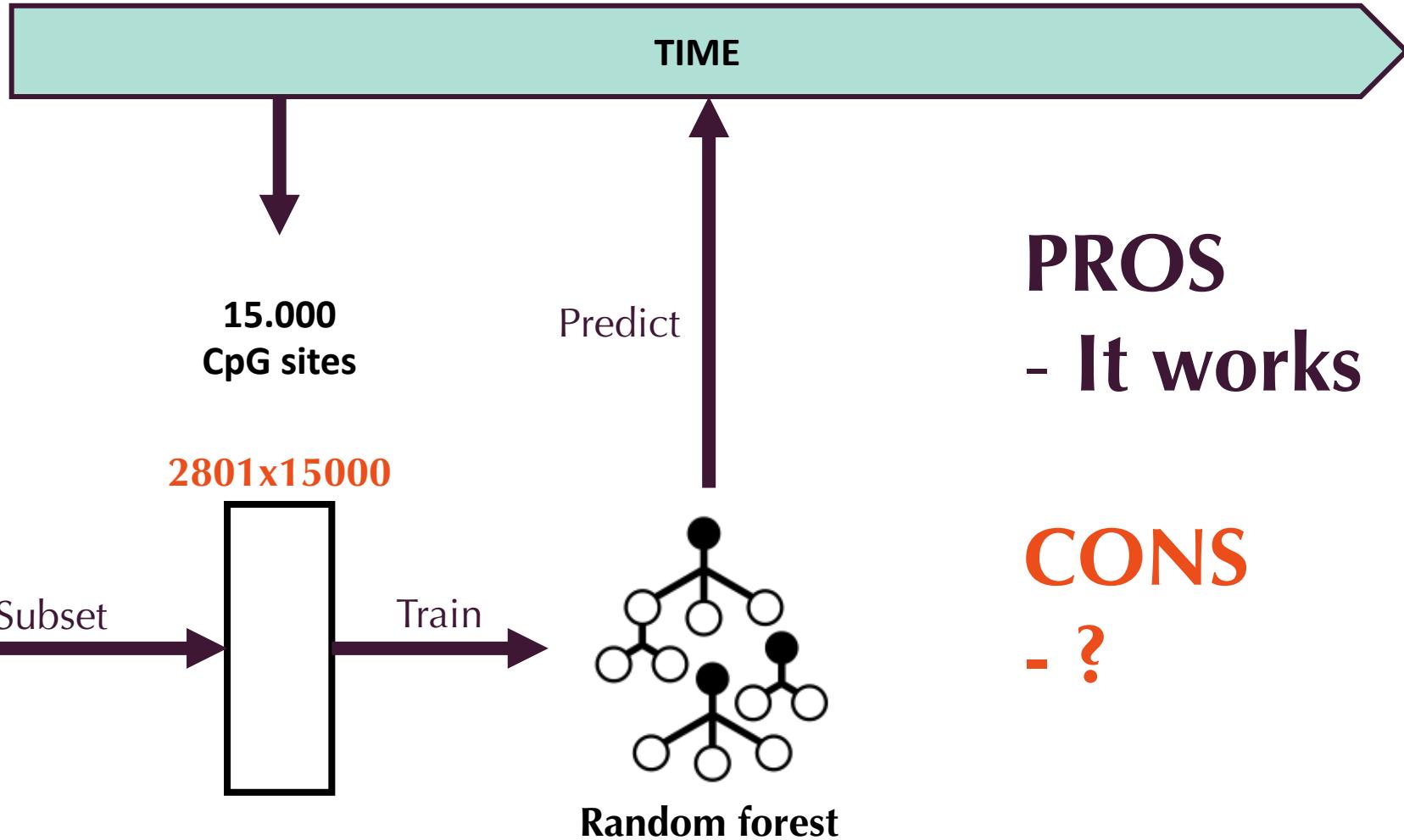
Existing approaches



Nanopore sequencing

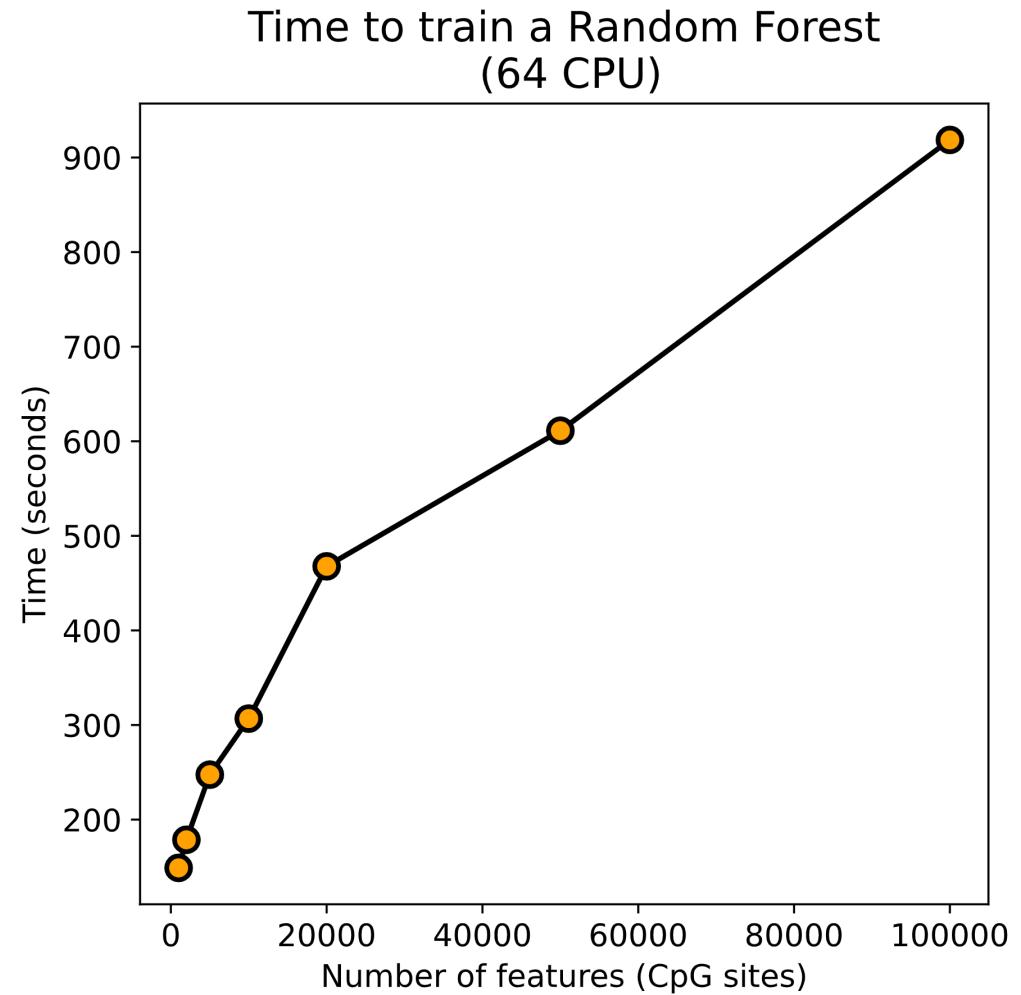


Microarray data



Existing approaches

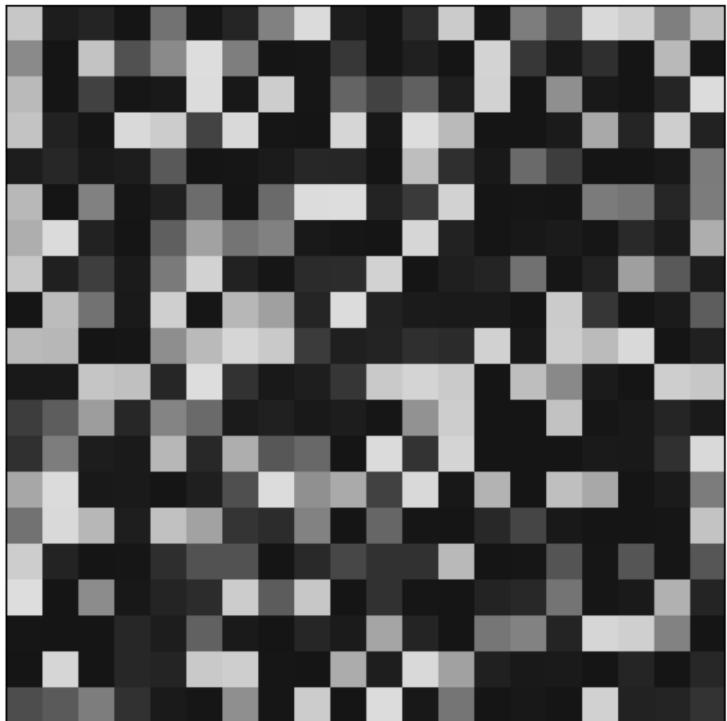
- Training takes valuable time
- Patient specific models
- Requires on-the-fly validation
- Requires access to data
- Poor scalability



Research question

Can we develop a model that can handle an arbitrary random amount of missing values?

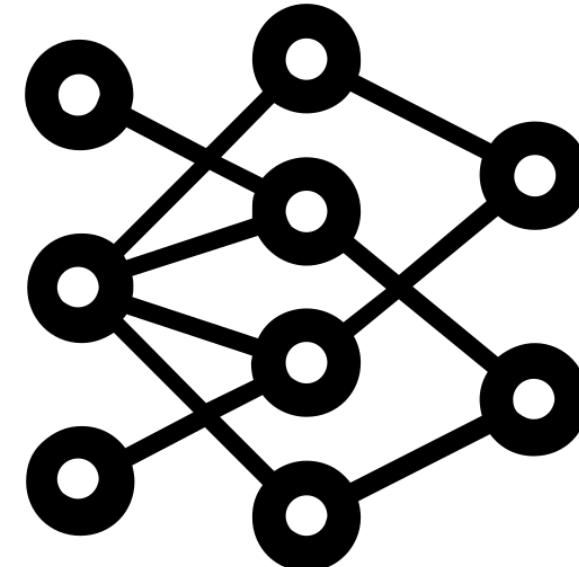
Our approach: Sturgeon



Micoarray data
2801 samples



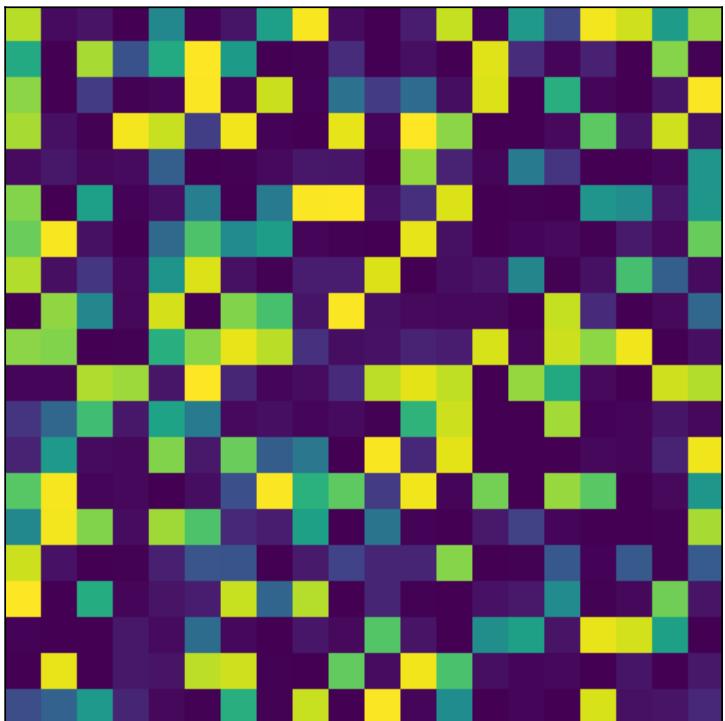
Simulate Nanopore data
> 36.000.000 simulations



Neural networks

Data simulation

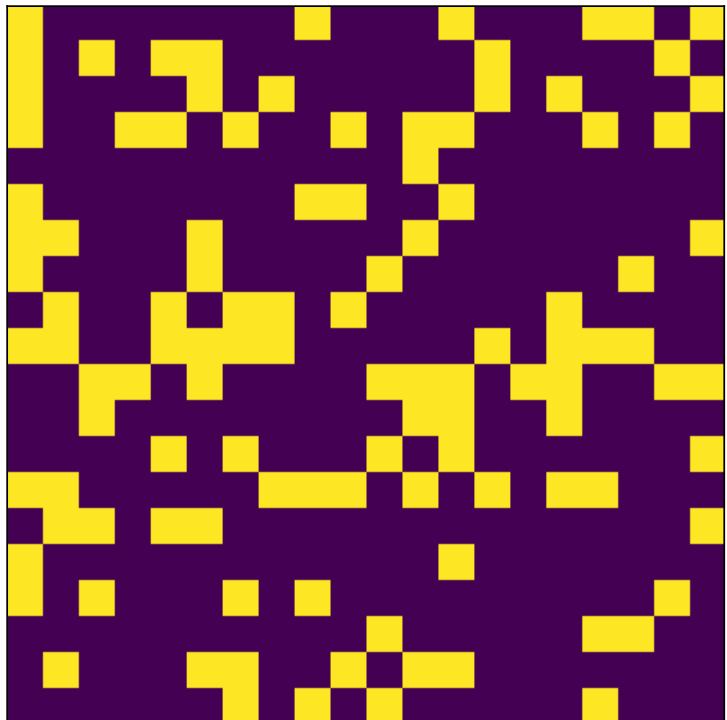
Microarray



Continuous measurement
Fraction of methylated reads

Data simulation

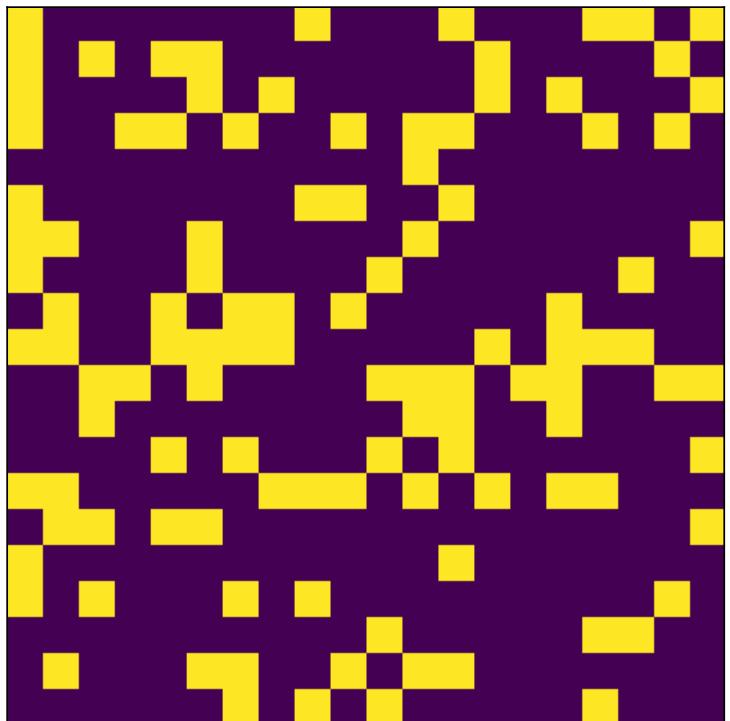
Microarray



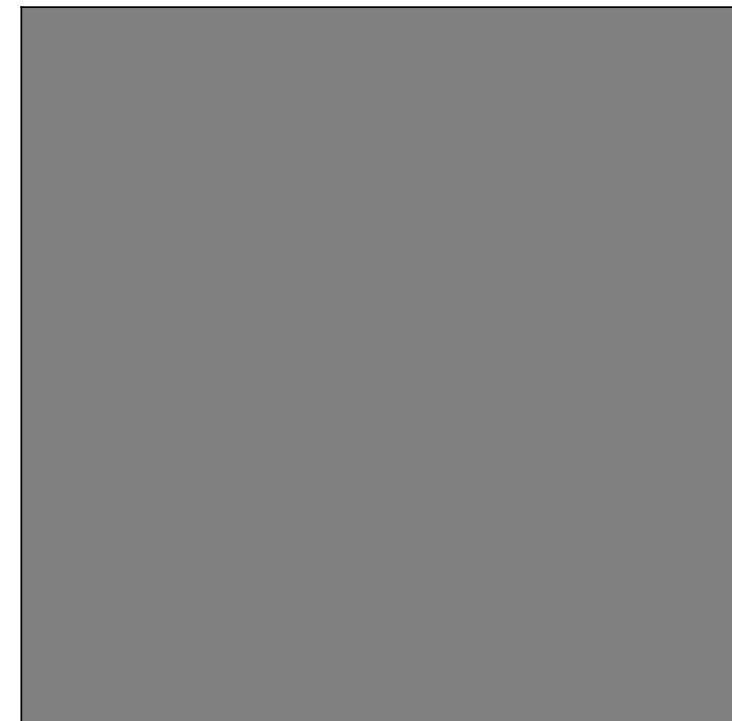
Binarized measurement
In shallow Nanopore sequencing
a site is methylated or not,
coverage of a site is very unlikely >1

Data simulation

Microarray

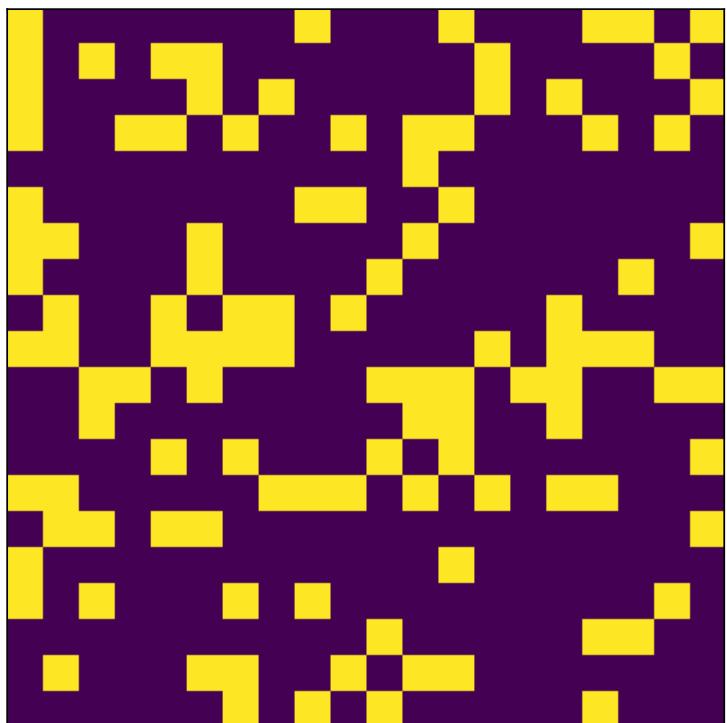


Nanopore



Data simulation

Microarray

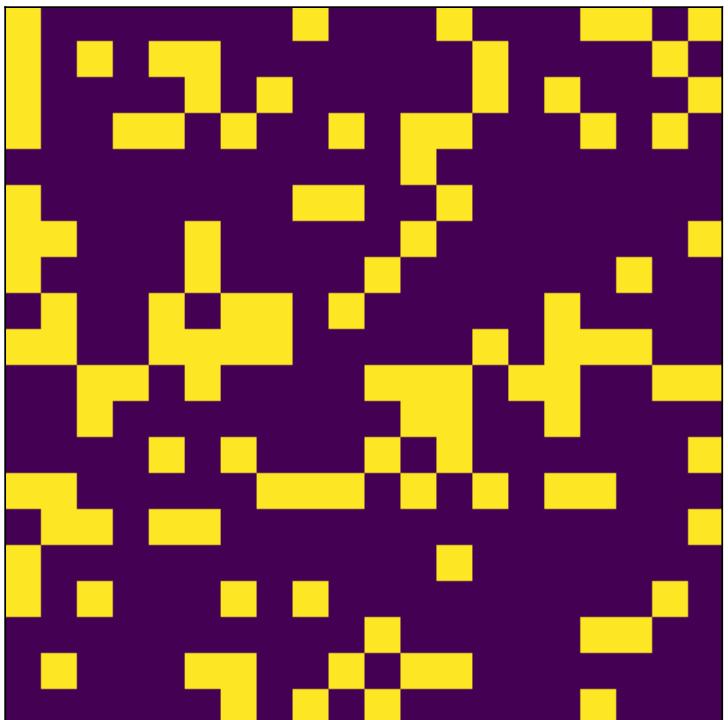


Nanopore



Data simulation

Microarray



Nanopore



Read simulation

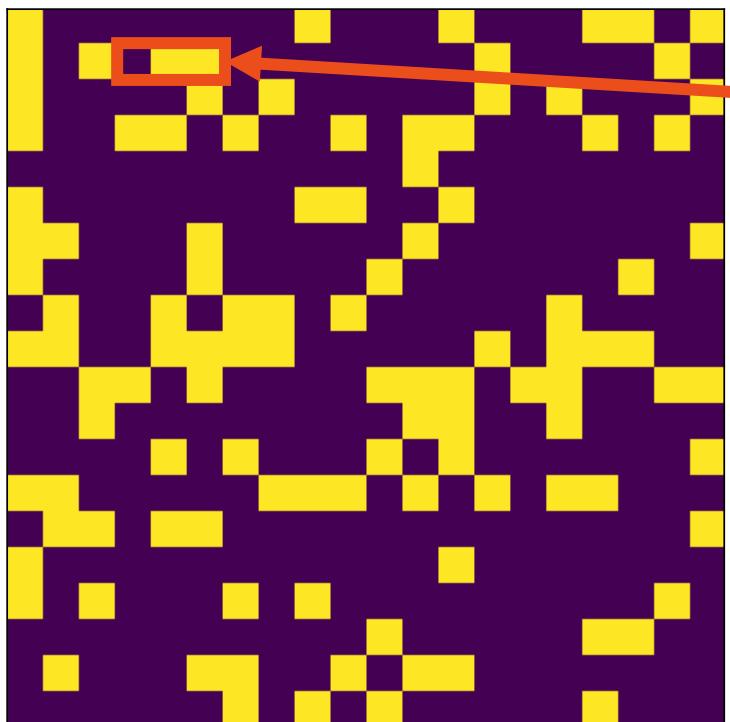
Genome position
[3 billion]

Strand
[Fwd, Rev]

Read length
[Distribution]

Data simulation

Microarray



Nanopore



Chr3
45000–46000
Fwd

x1

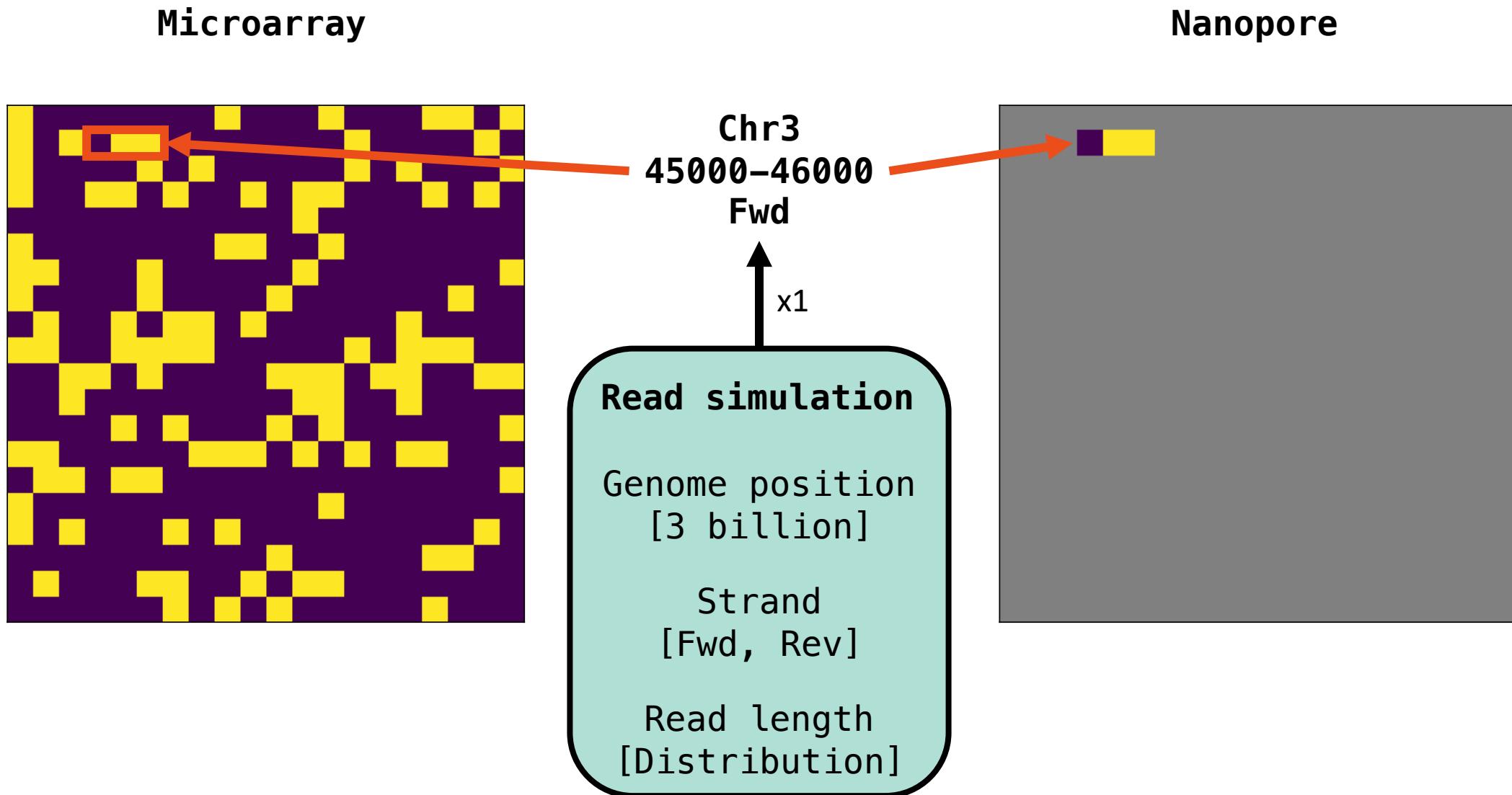
Read simulation

Genome position
[3 billion]

Strand
[Fwd, Rev]

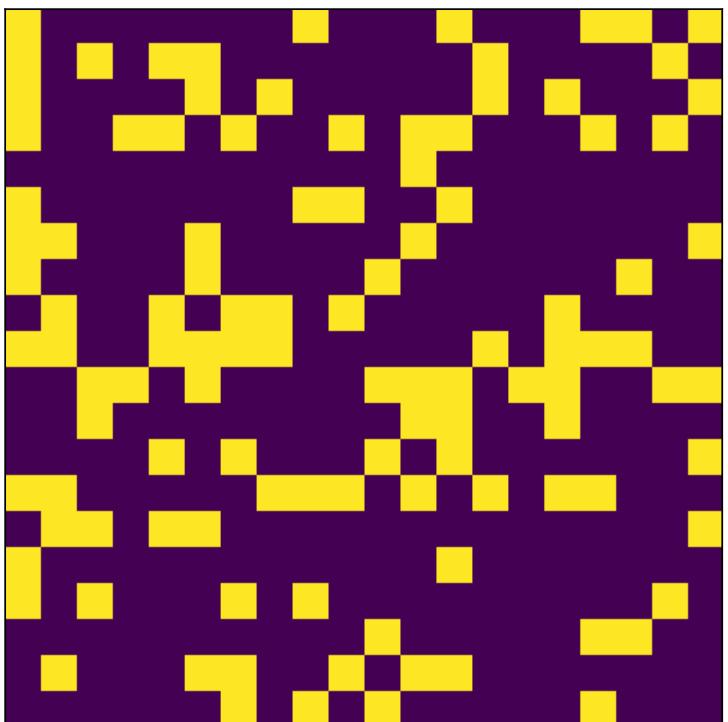
Read length
[Distribution]

Data simulation



Data simulation

Microarray



Nanopore



Chr17
30000–50000
Rev

x2

Read simulation

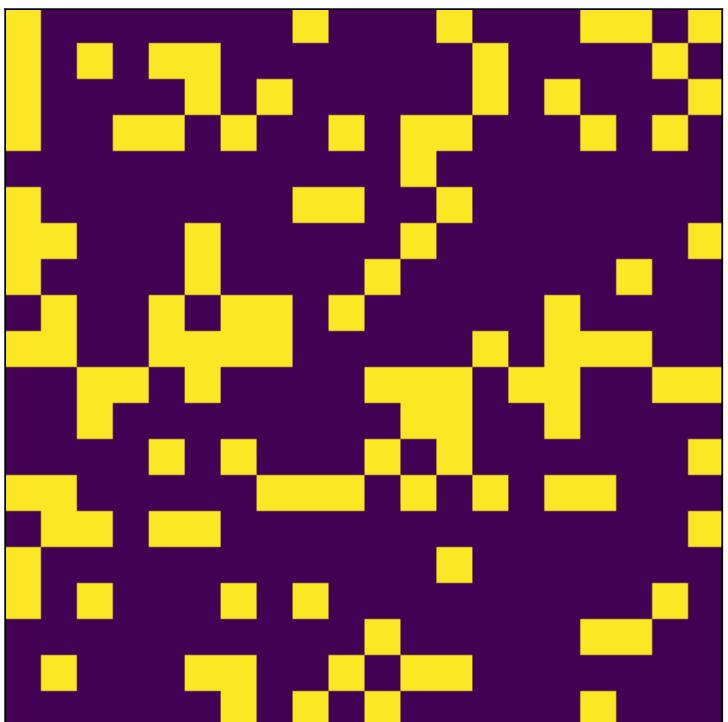
Genome position
[3 billion]

Strand
[Fwd, Rev]

Read length
[Distribution]

Data simulation

Microarray



Nanopore



Chr18
2000–10000
Rev

x3

Read simulation

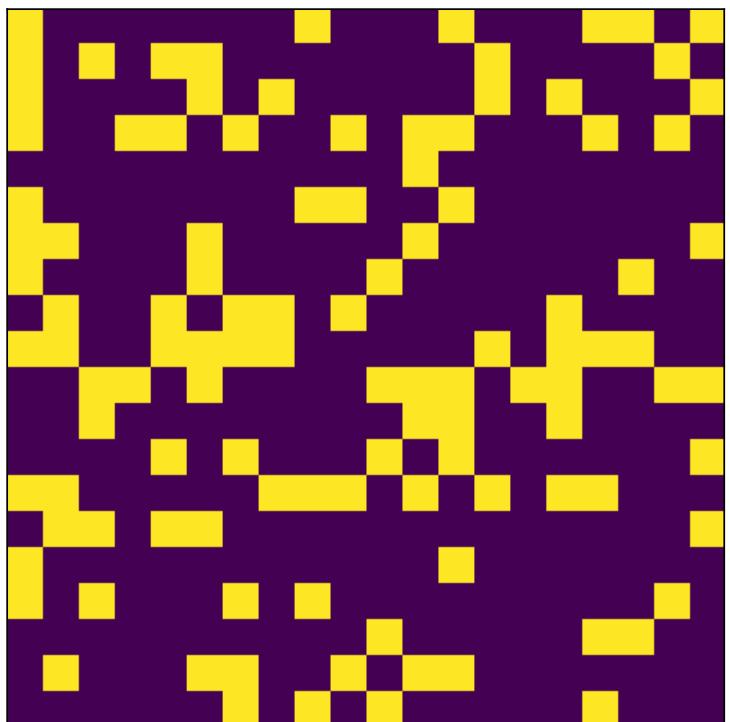
Genome position
[3 billion]

Strand
[Fwd, Rev]

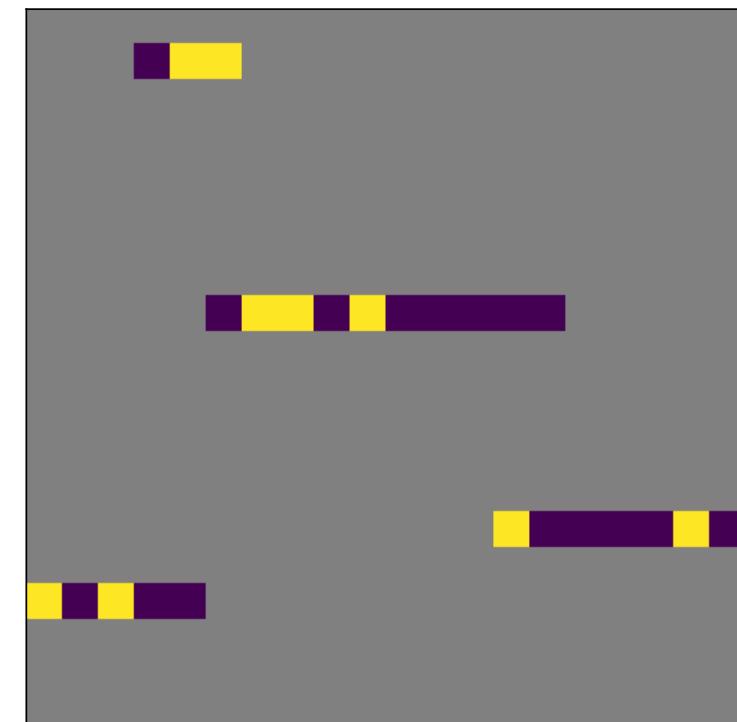
Read length
[Distribution]

Data simulation

Microarray



Nanopore



Chr10
12000–20000
Rev

x4

Read simulation

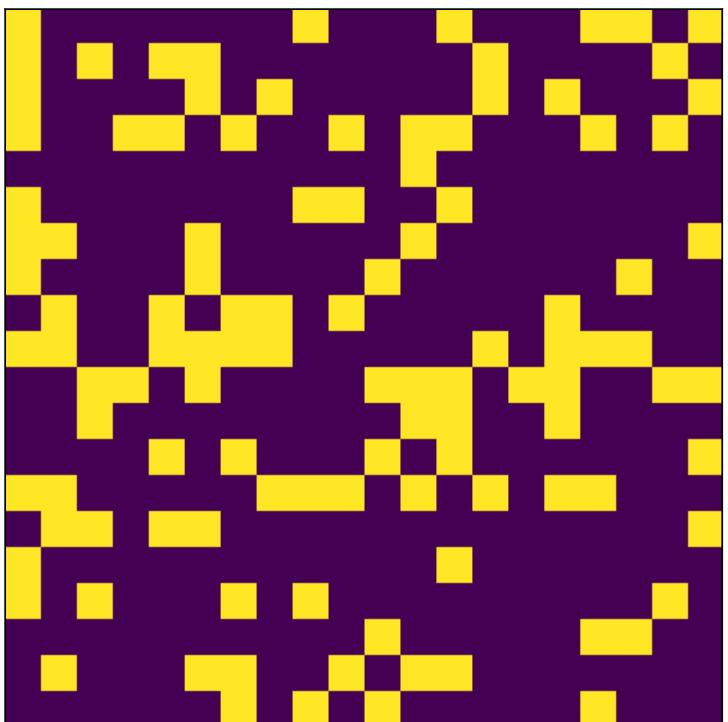
Genome position
[3 billion]

Strand
[Fwd, Rev]

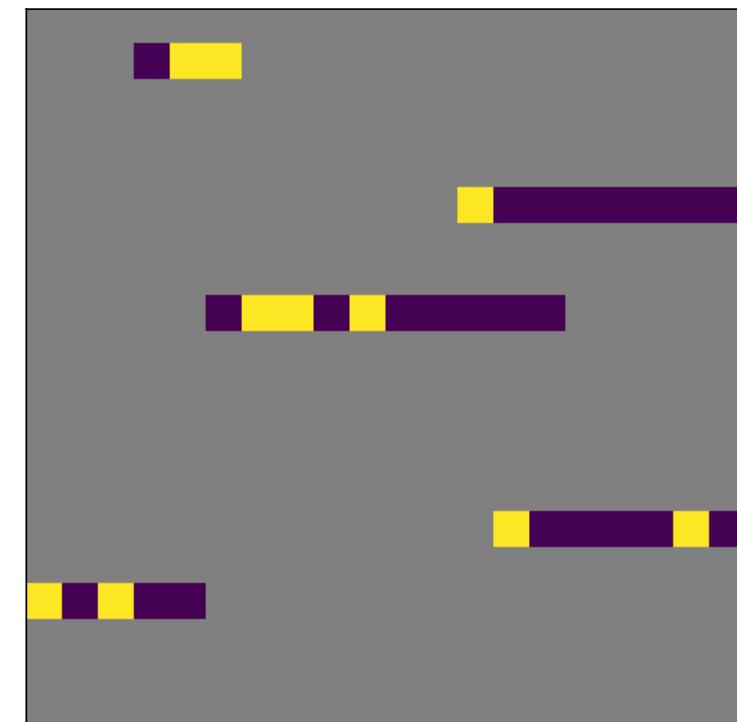
Read length
[Distribution]

Data simulation

Microarray



Nanopore



Chr8
25000–45000
Fwd

x5

Read simulation

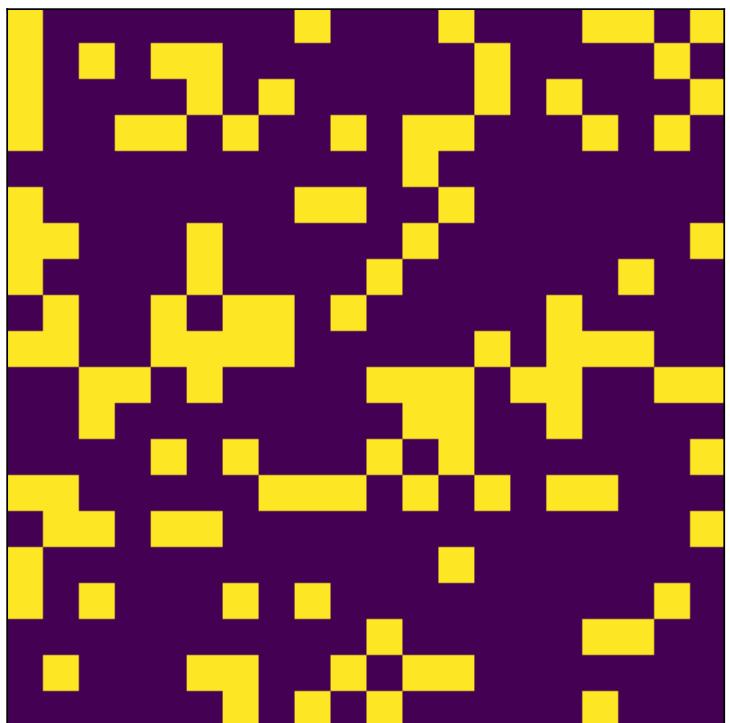
Genome position
[3 billion]

Strand
[Fwd, Rev]

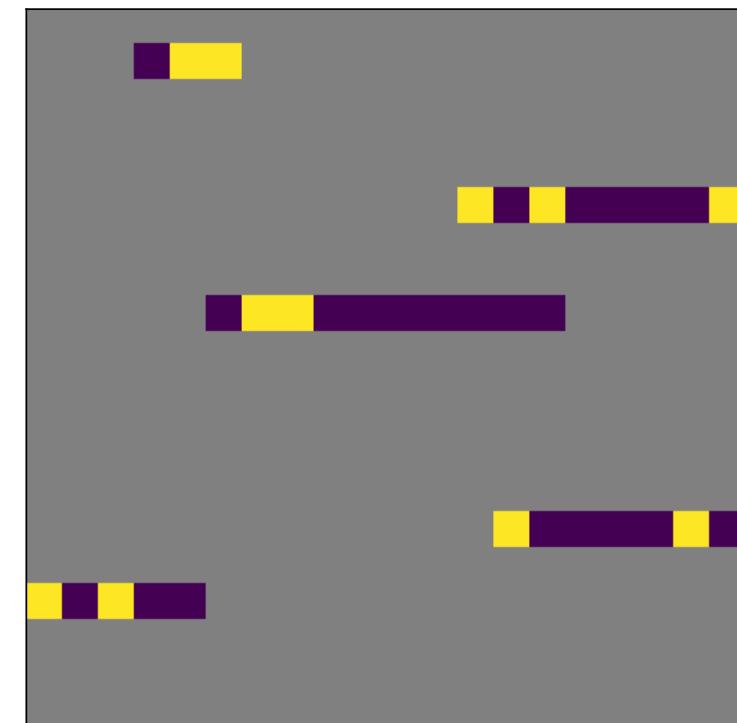
Read length
[Distribution]

Data simulation

Microarray



Nanopore



Chr8
25000–45000
Fwd

x5

Read simulation

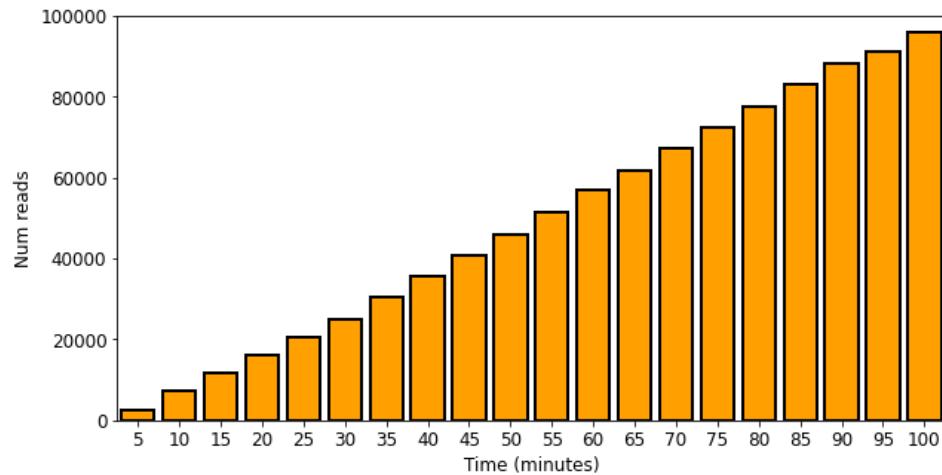
Genome position
[3 billion]

Strand
[Fwd, Rev]

Read length
[Distribution]

+10% noise
Flip methylation status

Data simulation



How many reads do we simulate?

How much time,
do we want to sequence?

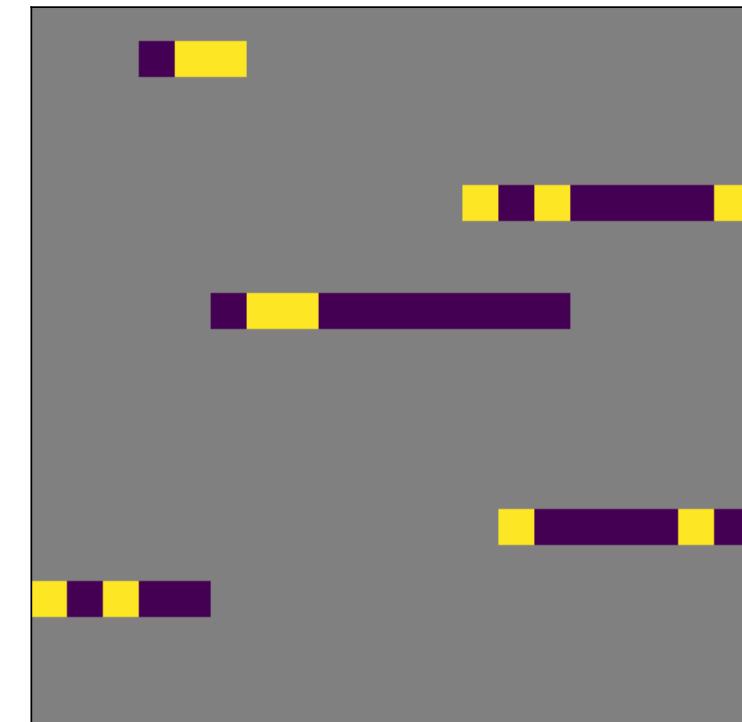
ChrN
Start-End
Strand

Read simulation

Genome position
[3 billion]

Strand
[Fwd, Rev]

Read length
[Distribution]

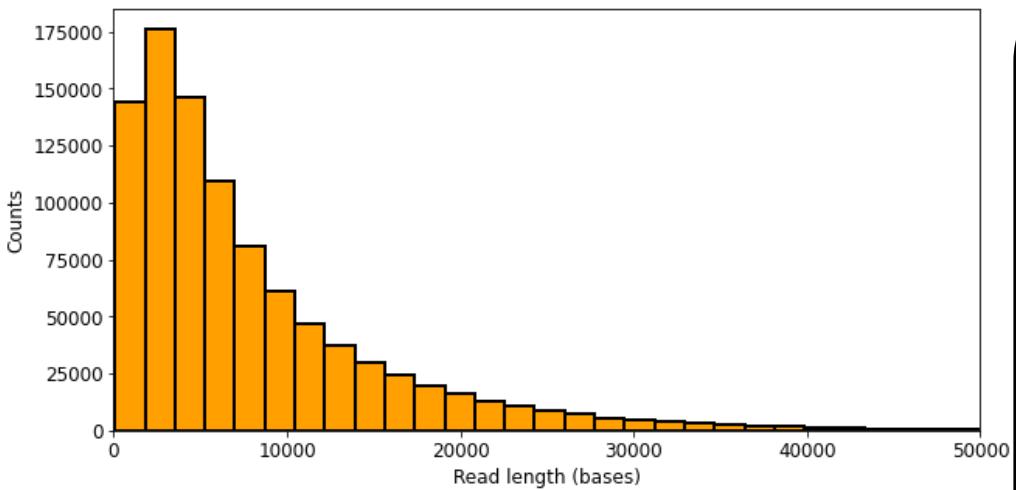


Nanopore

+10% noise
Flip methylation status

Data simulation

How long should be the simulated reads?



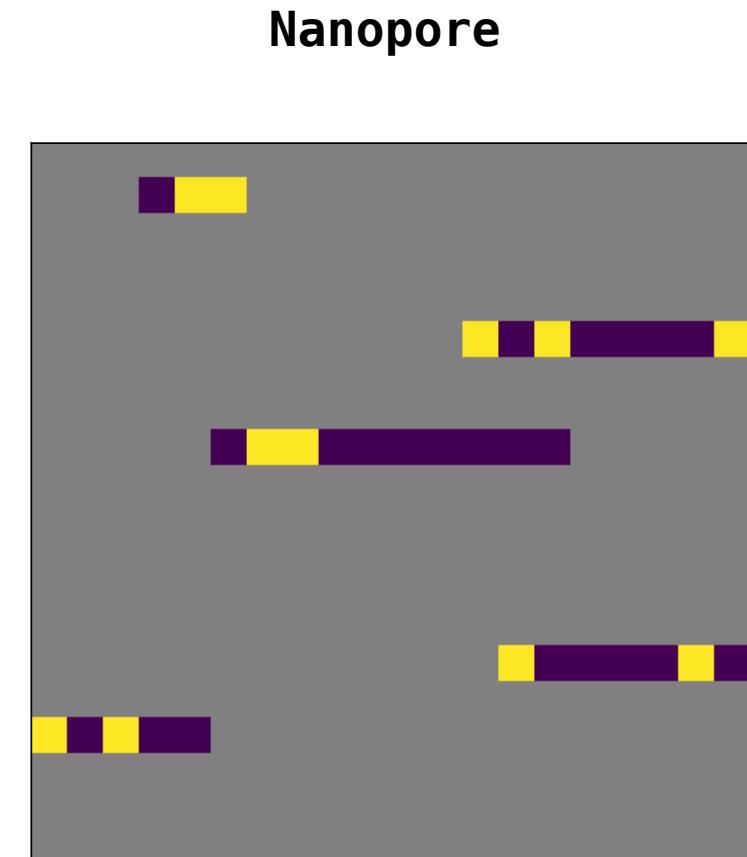
ChrN
Start-End
Strand
xN

Read simulation

Genome position
[3 billion]

Strand
[Fwd, Rev]

Read length
[Distribution]



+10% noise
Flip methylation status

Neural network

426.000 input channels

Linear layer -> Sigmoid

256 hidden channels

Linear layer -> Sigmoid

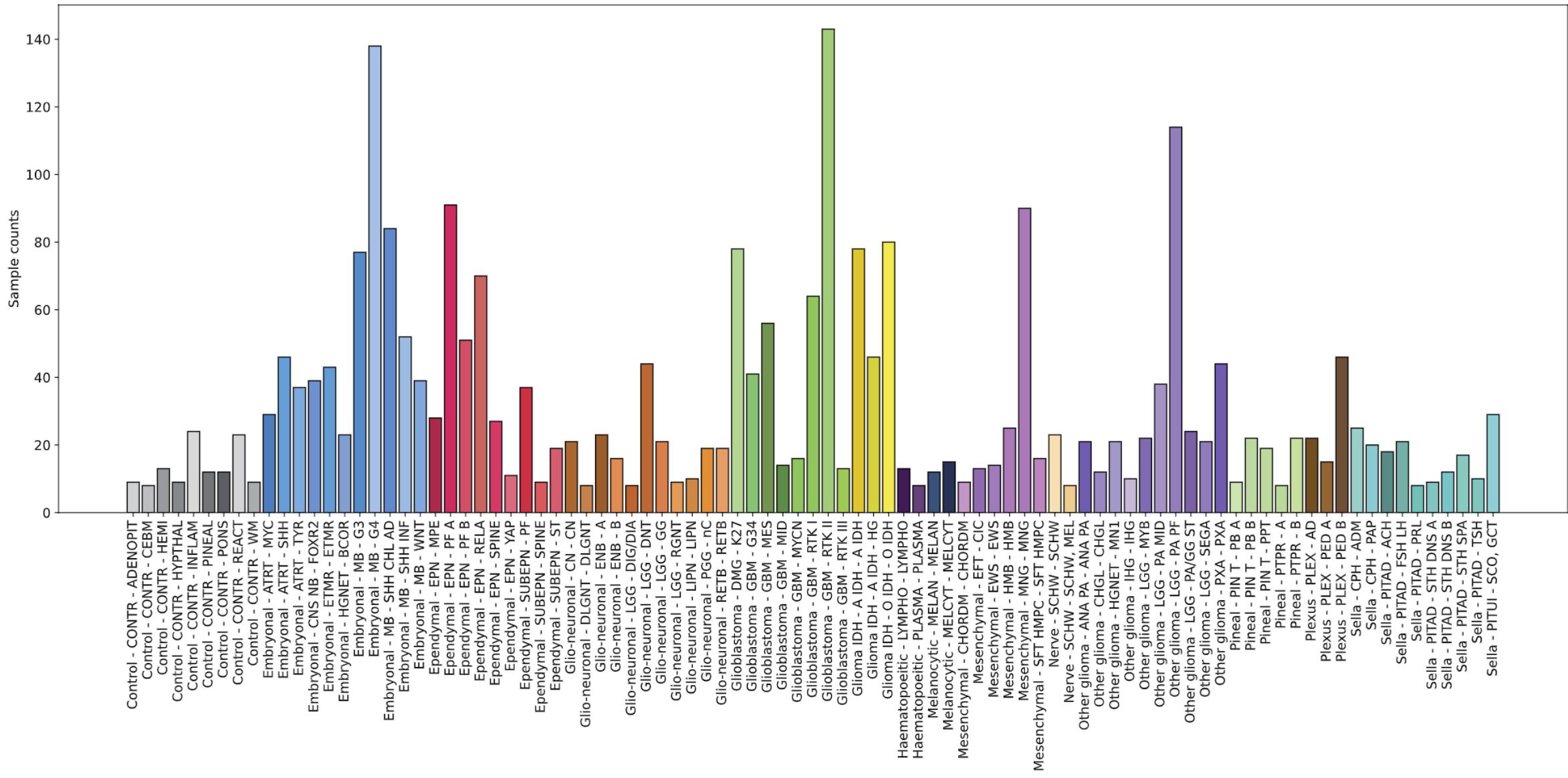
128 hidden channels

Linear layer

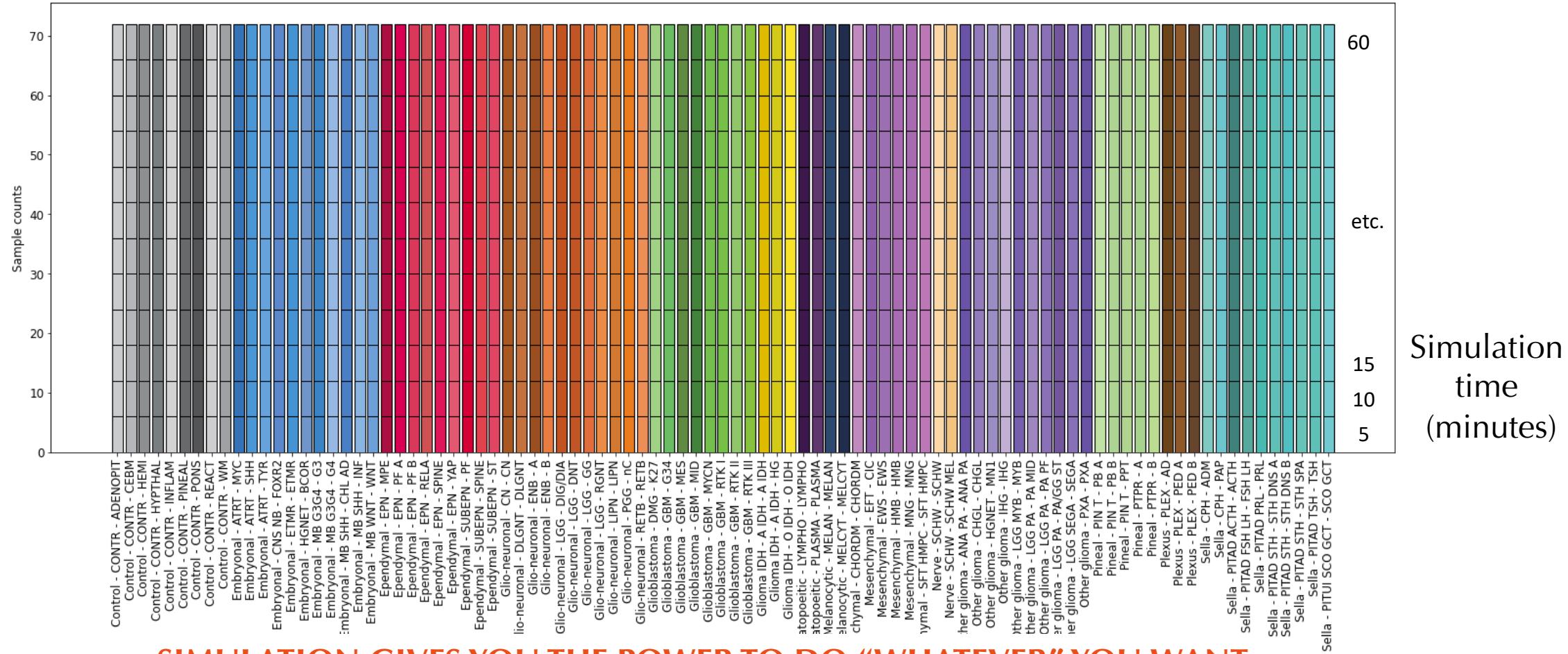
91 output channels

**Extremely simple feed
forward network
because the data has no
structure**

Cross-validation



Class balancing via up-sampling

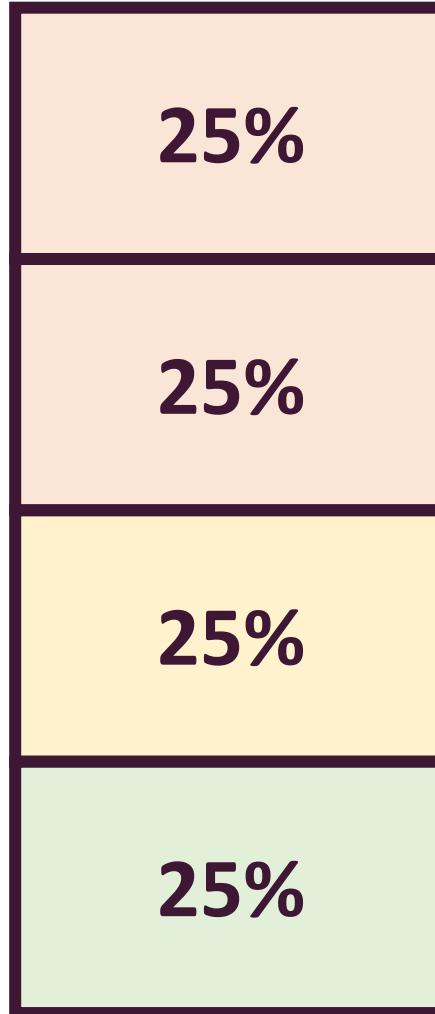
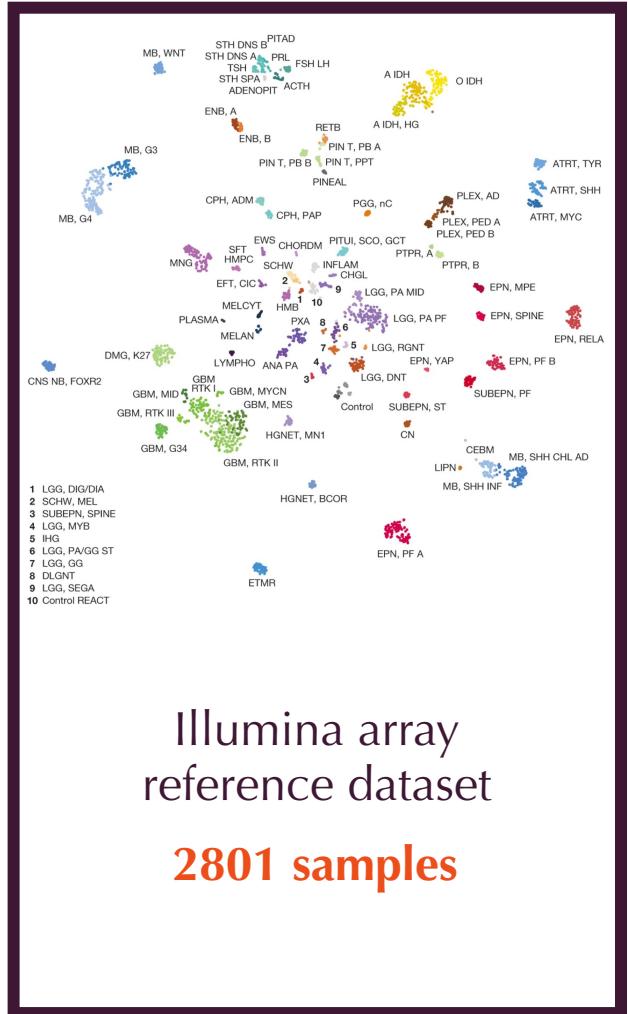


SIMULATION GIVES YOU THE POWER TO DO “WHATEVER” YOU WANT

Upsample all classes to the number of samples of the most represented class
 Equal simulation times for each class and number of samples

Cross-validation

Keep balanced class proportions

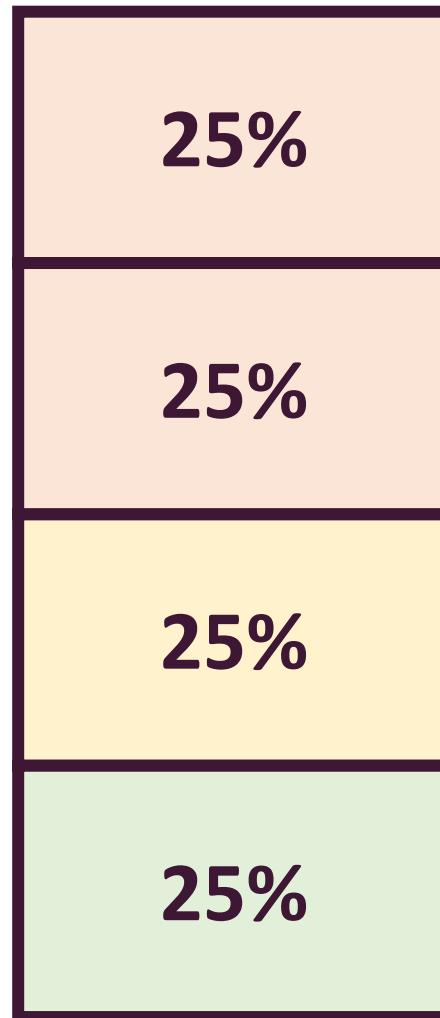
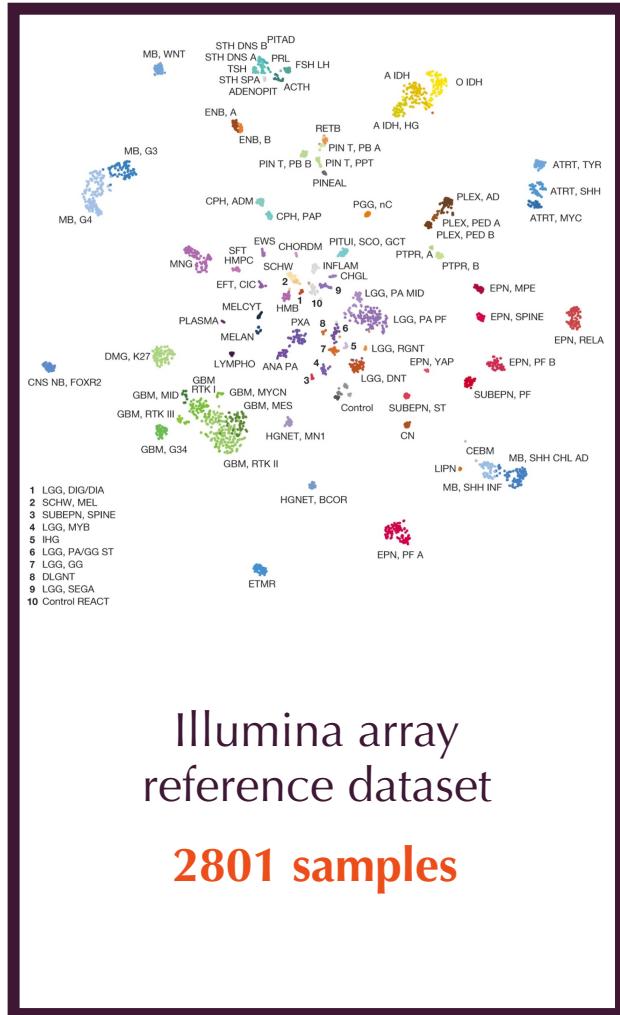


Stop training
Model calibration

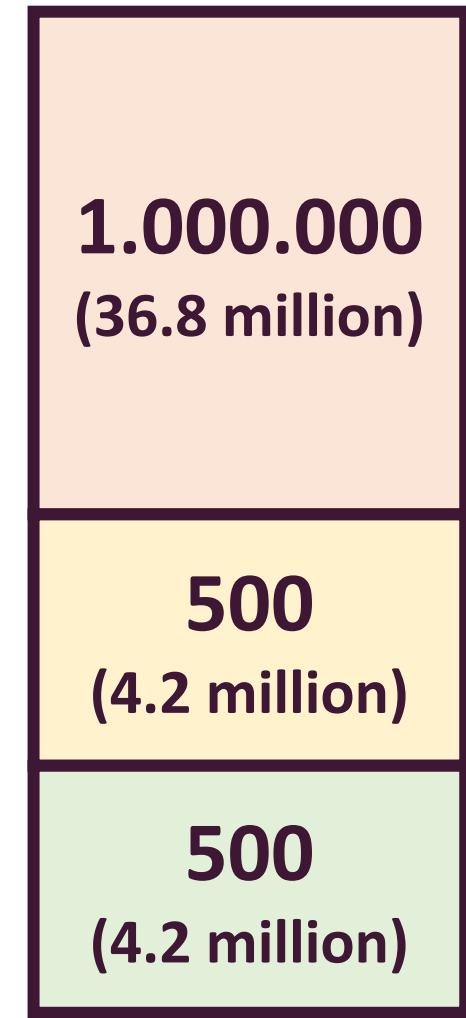
Final model performance

Cross-validation

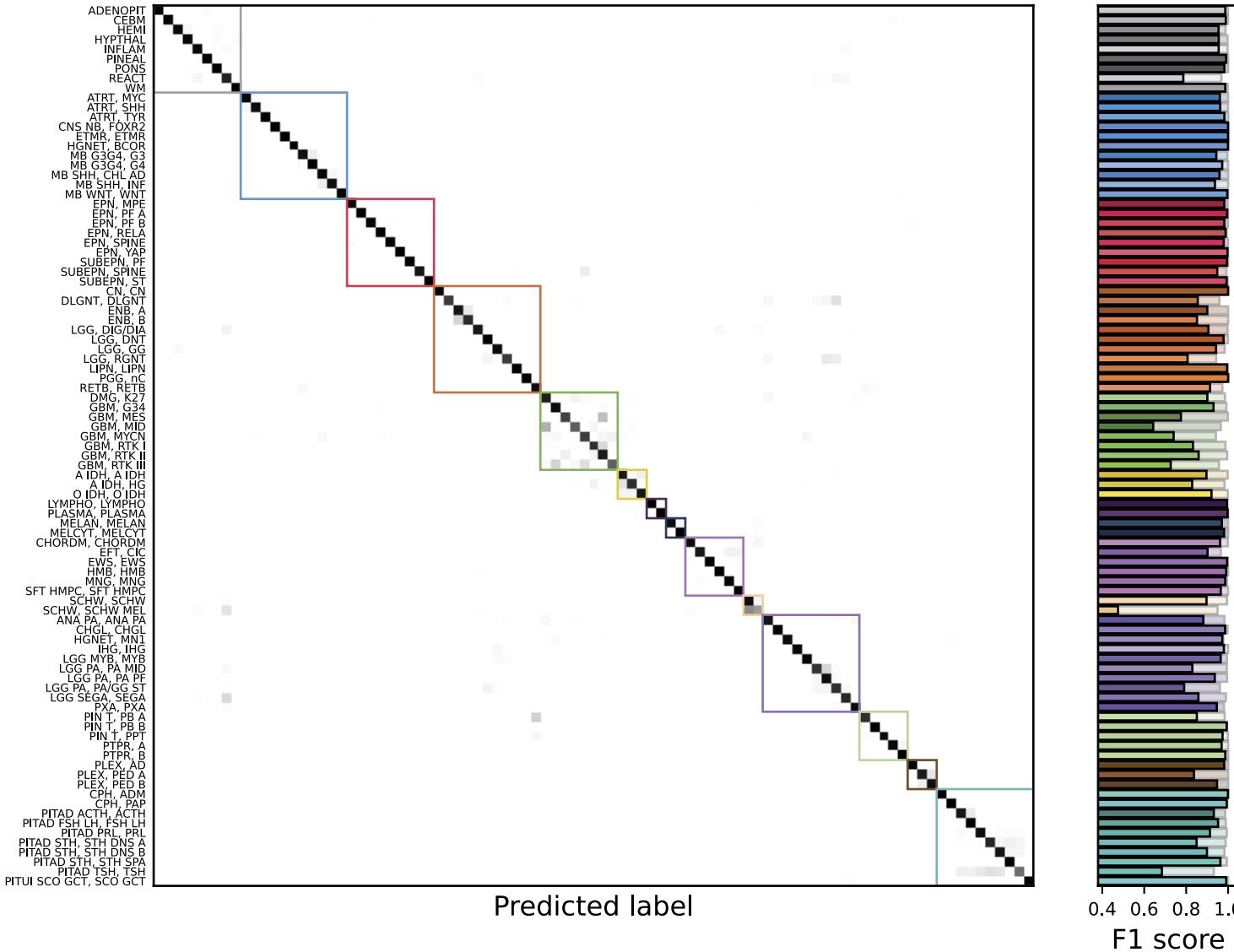
Keep balanced class proportions



Simulation seeds (non-overlapping)



Performance on simulations

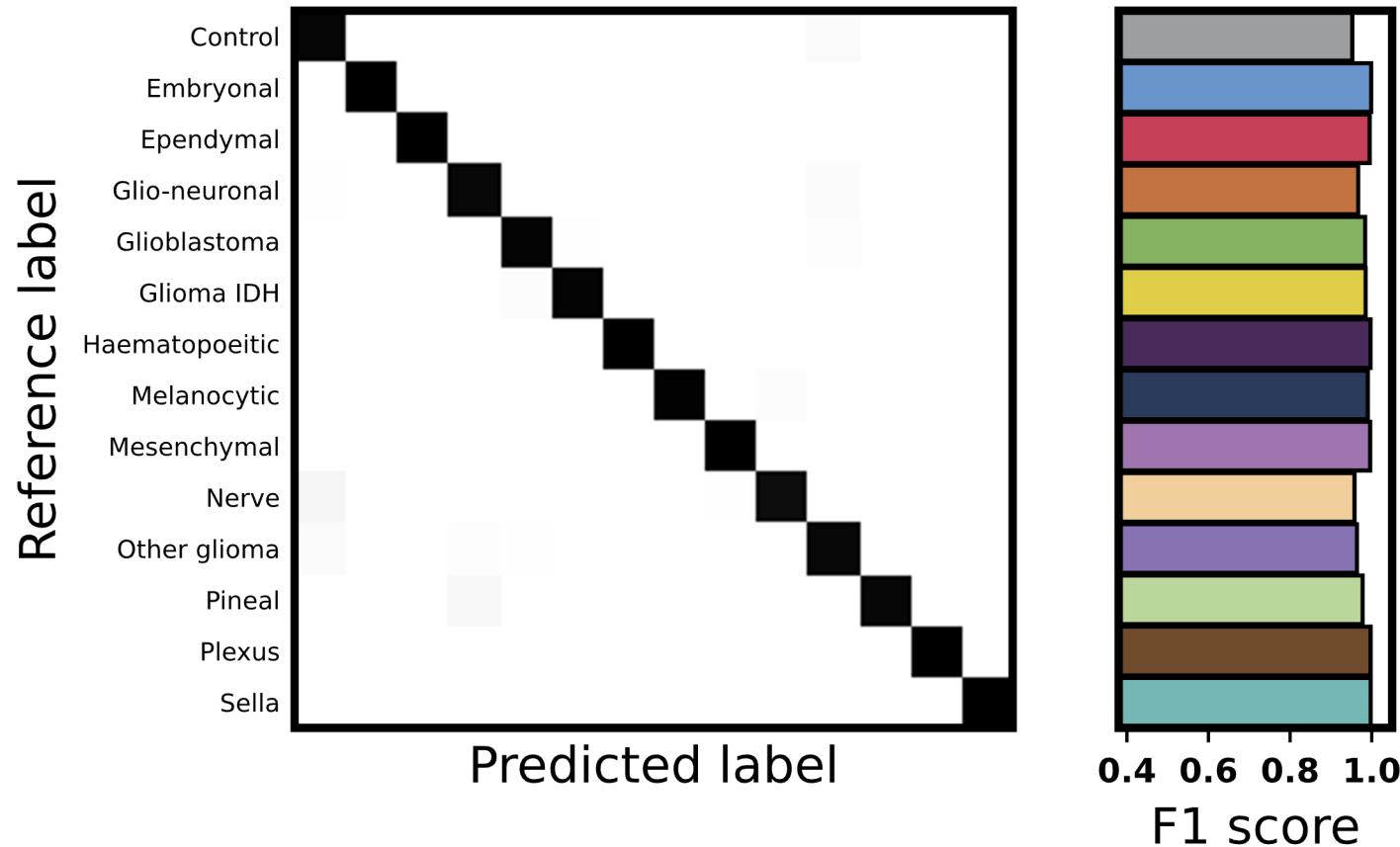


Performance at **20 minutes** simulations

F1-Score
Top1: 0.926
Top3: 0.990

$$F1 = \frac{2TP}{2TP + FP + FN}$$

Performance on simulations

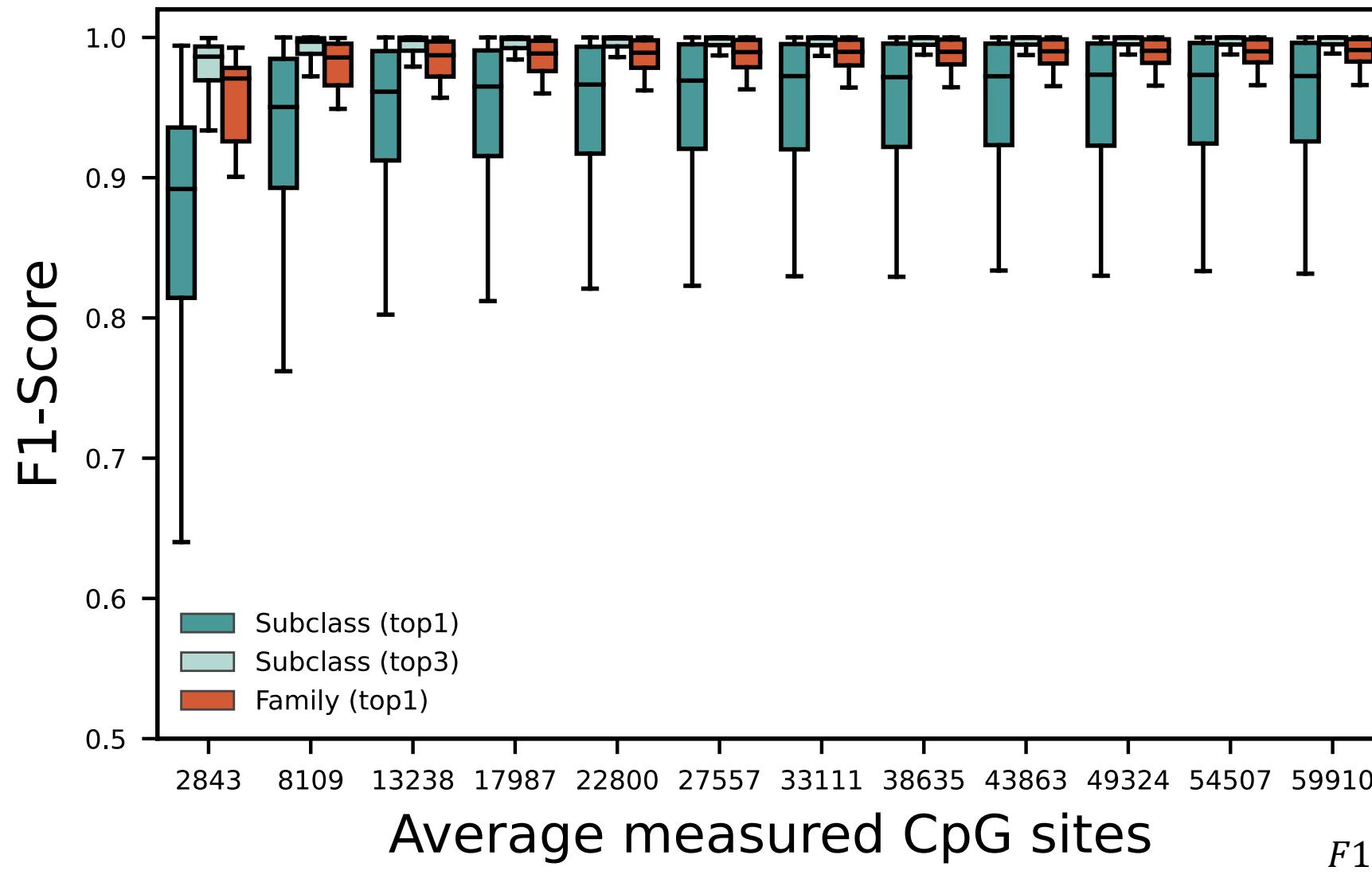


Performance at **20 minutes** simulations

F1-Score
Top1: 0.982

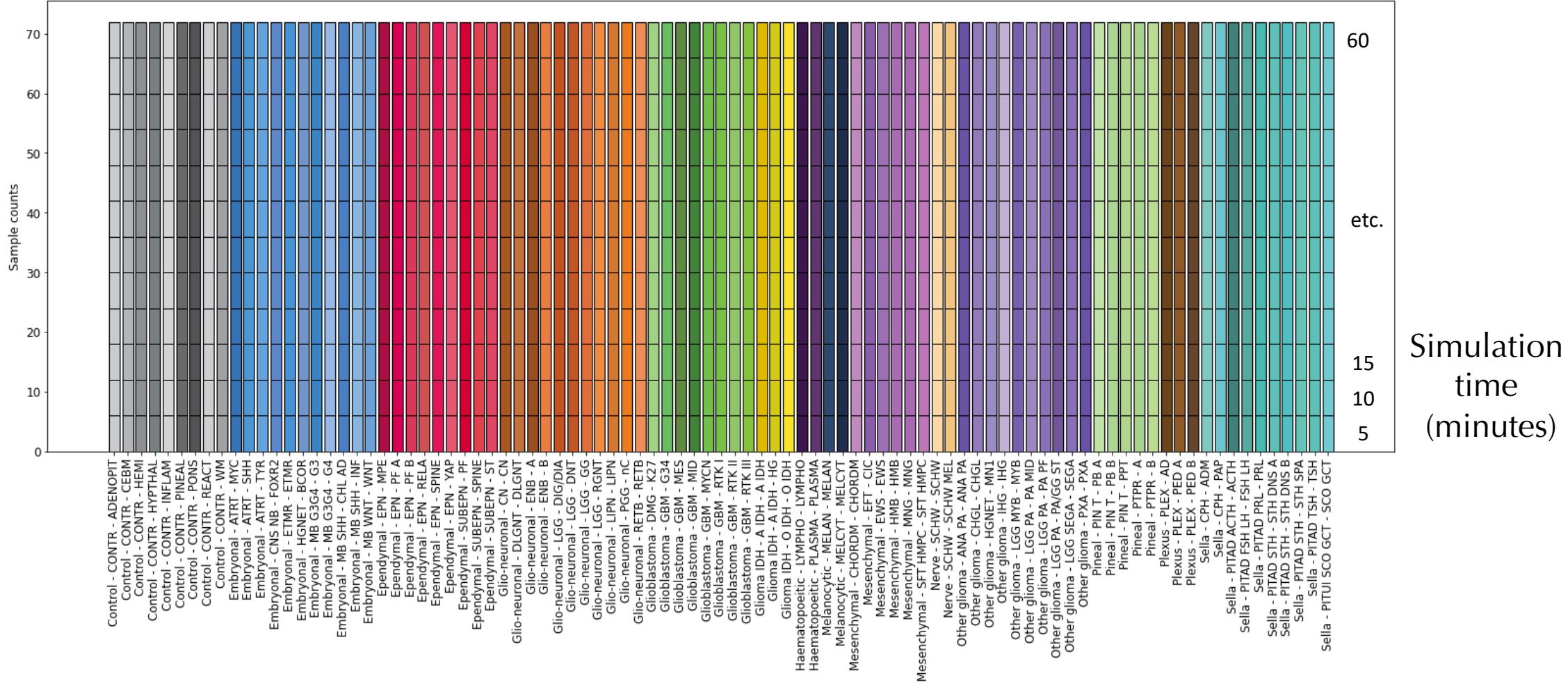
$$F1 = \frac{2TP}{2TP + FP + FN}$$

Performance on simulations



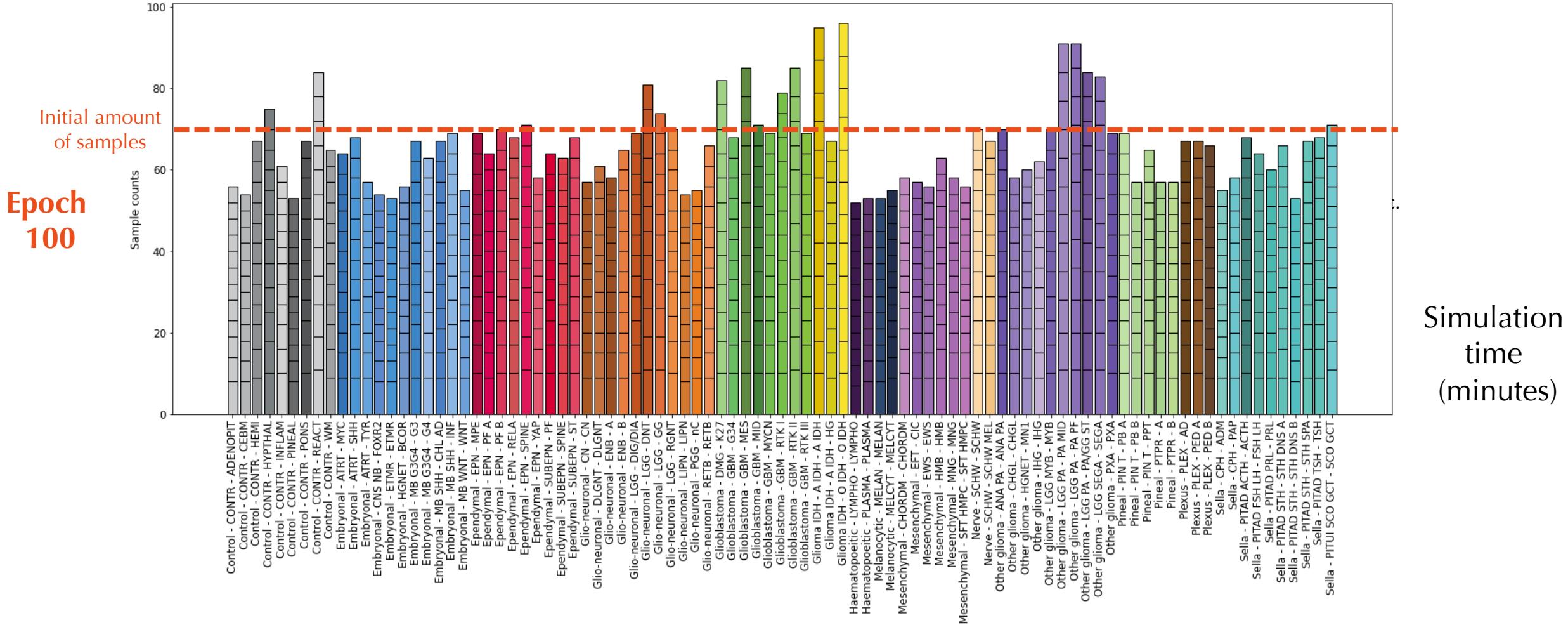
Adaptive sampling

Epoch
1



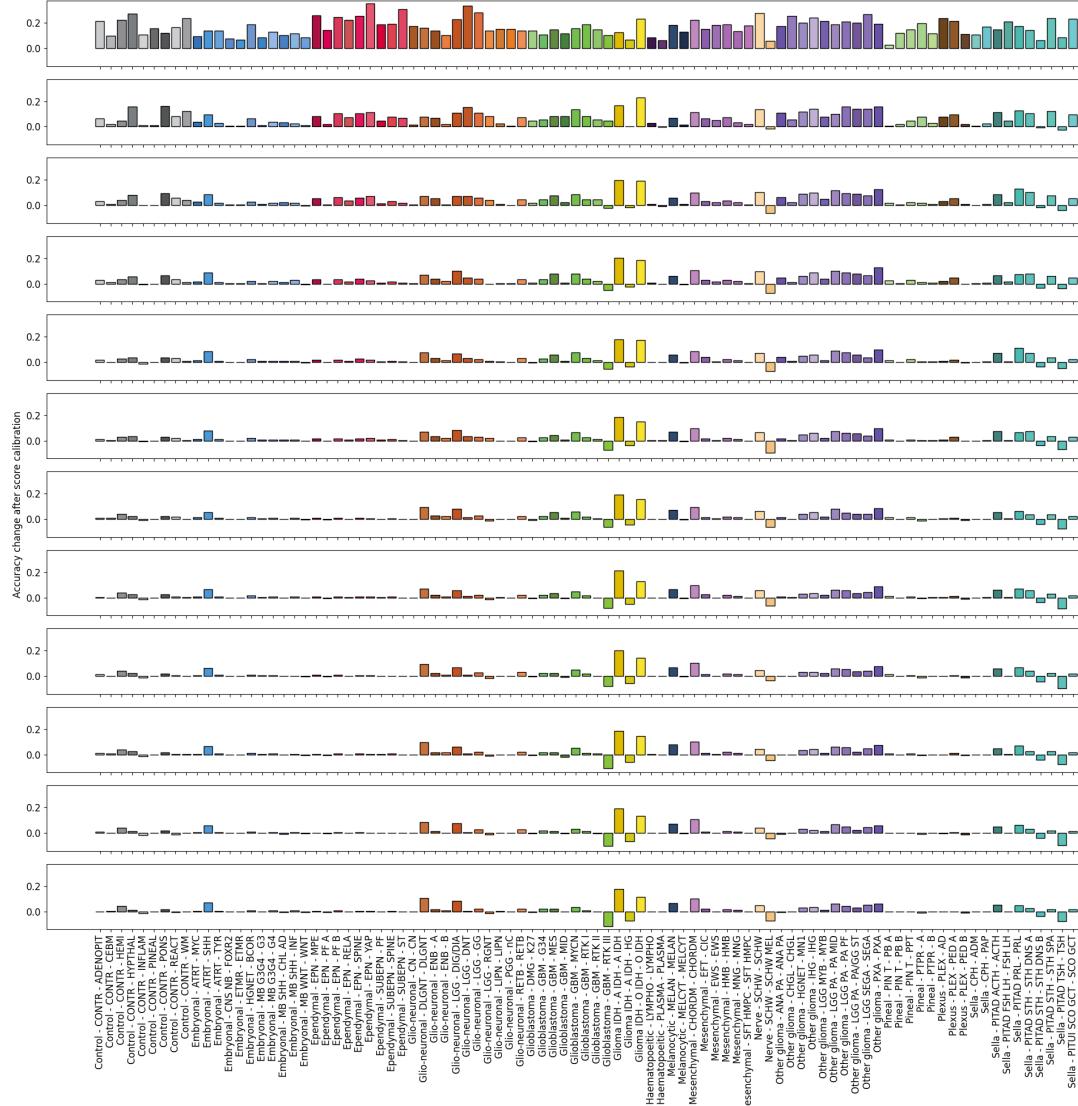
Give homework on the classes and simulation times that the model does worse

Adaptive sampling



Give homework on the classes and simulation times that the model does worse

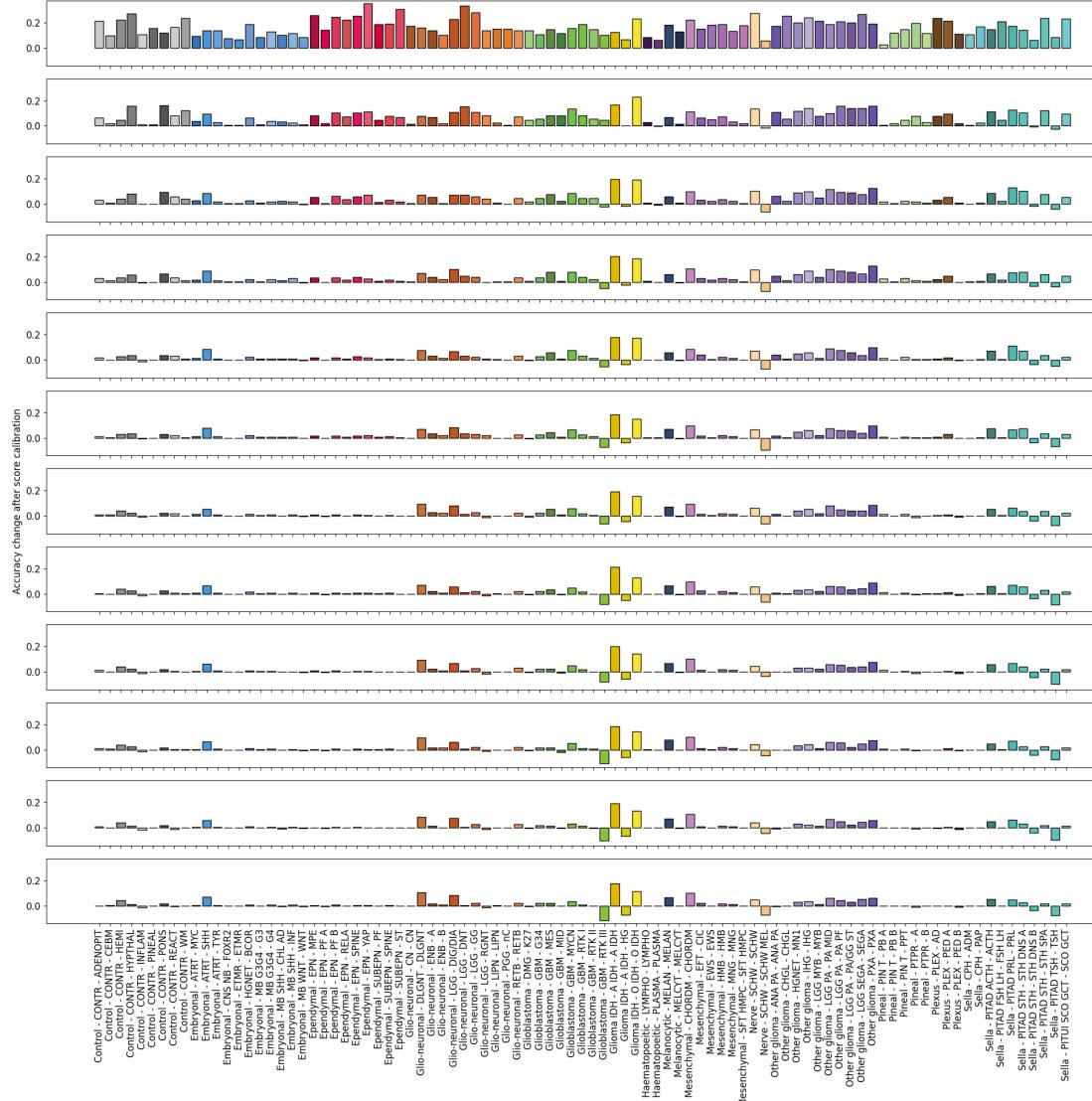
Adaptive sampling



5 min

Adaptive sampling
improves performance
at lower simulation times

Adaptive sampling

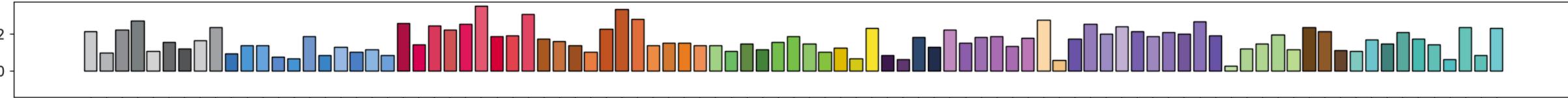


5 min

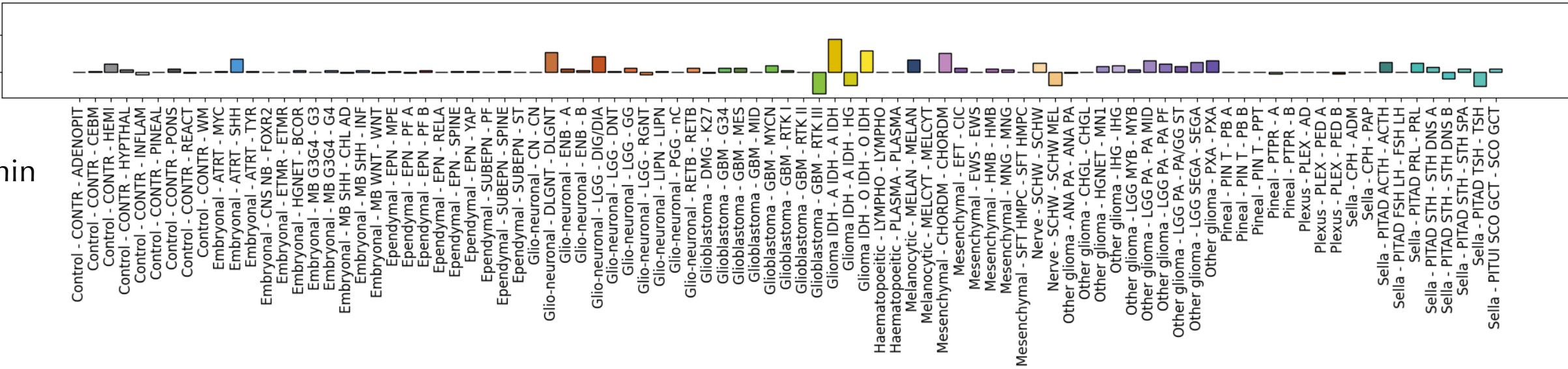
60 min

Adaptive sampling
improves performance
at lower simulation times

Adaptive sampling



5 min

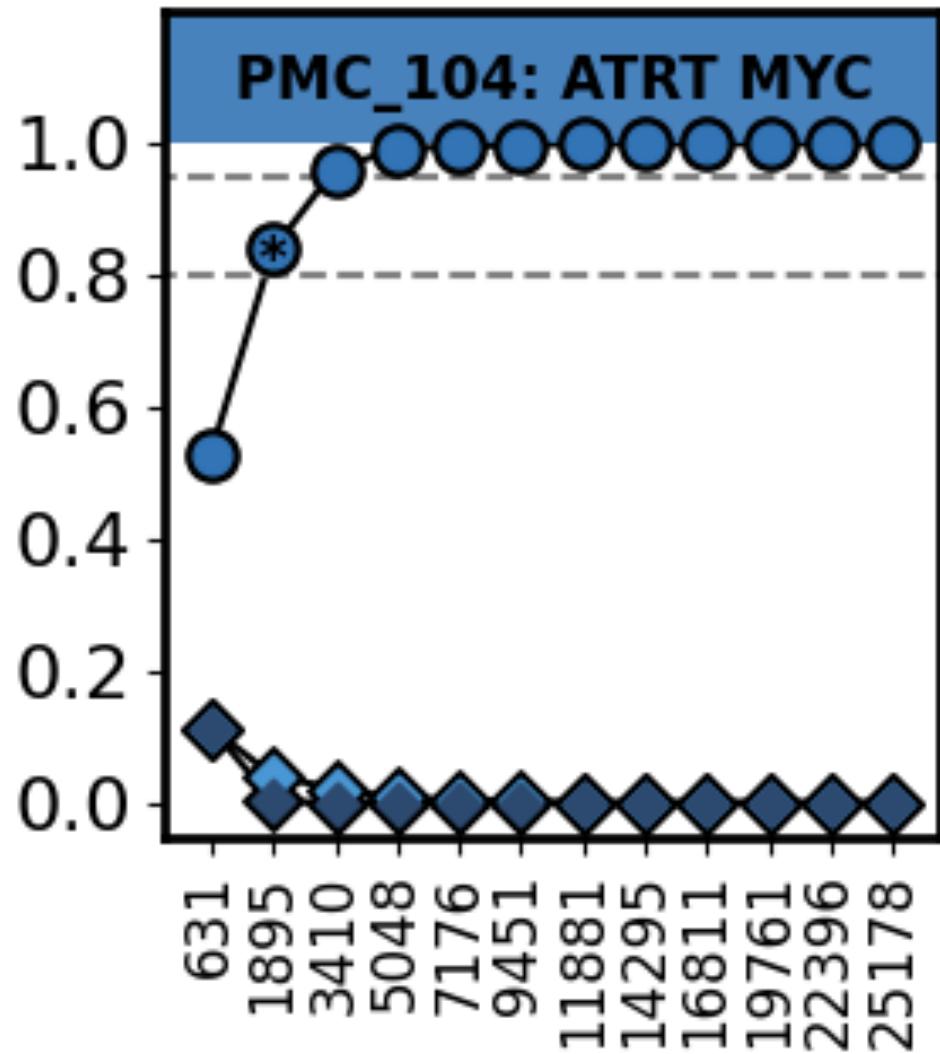


60 min

Nanopore

All this is nice, but it's all simulations
Does it actually work on Nanopore data

PMC Nanopore

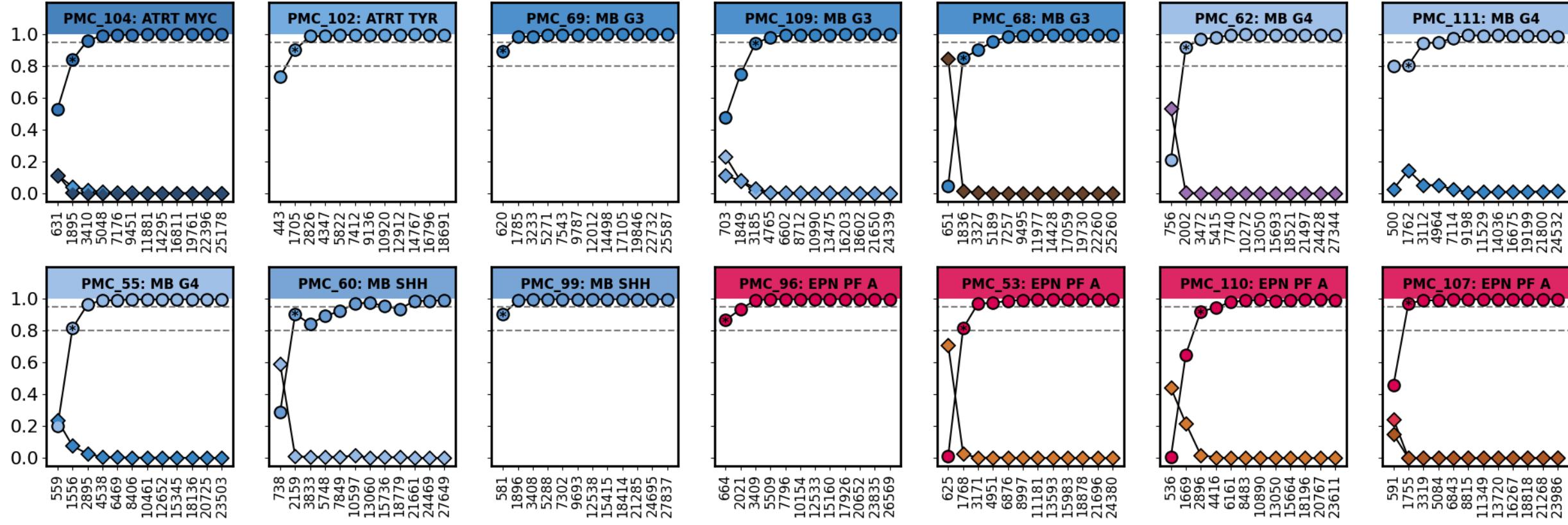


Circles = correct class

Rombes = incorrect class

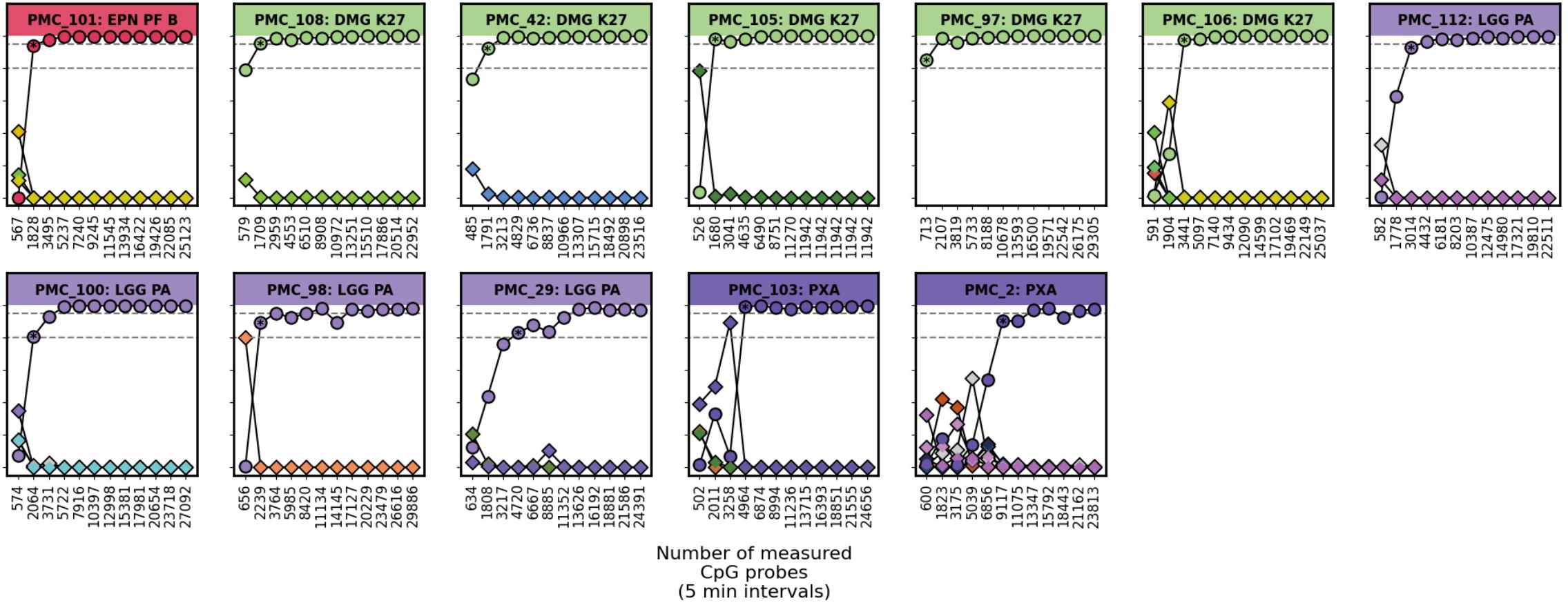
Each tick is 5 minutes of sequencing

PMC Nanopore

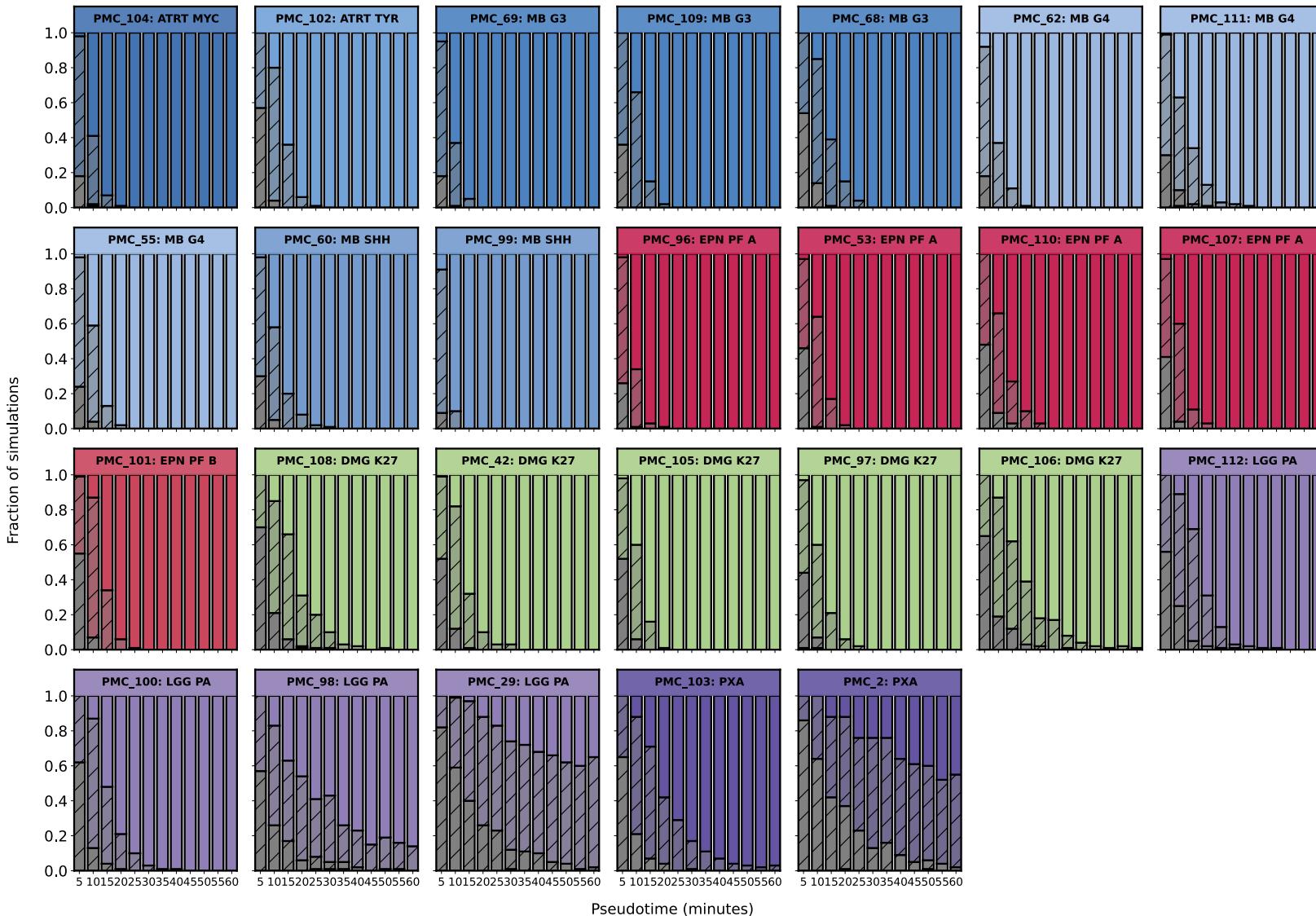


PMC Nanopore

Predicted s_t



Simulation power



Oslo Nanopore

ORIGINAL ARTICLE

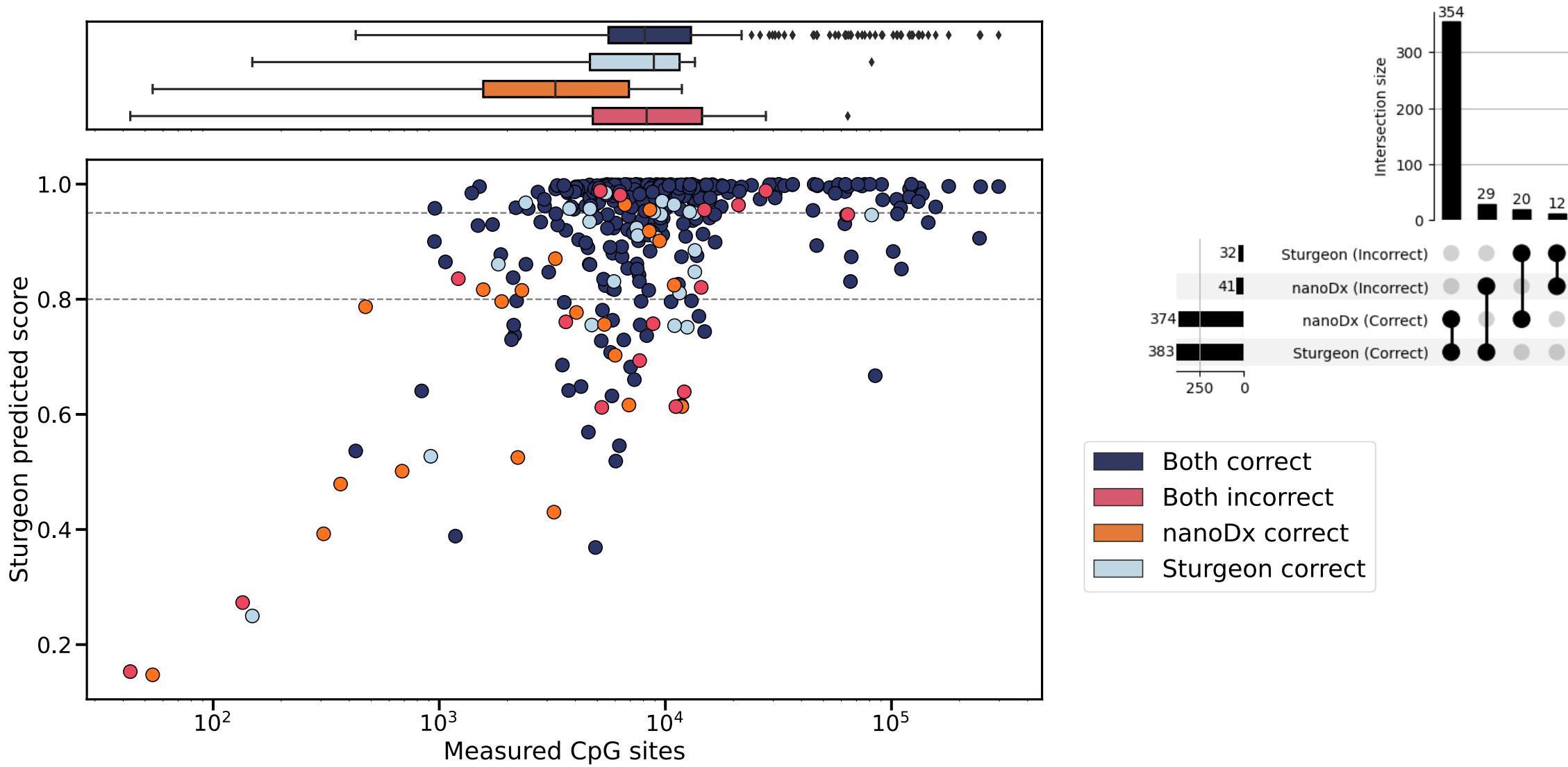
Neuropathology and
Applied Neurobiology
JOURNAL OF THE BRITISH NEUROPATHOLOGICAL SOCIETY

WILEY

Robust methylation-based classification of brain tumours using nanopore sequencing

Luis P. Kuschel¹  | Jürgen Hench² | Stephan Frank² | Ivana Bratic Hench² |
Elodie Girard³ | Maud Blanluet³ | Julien Masliah-Planchon³ | Martin Misch⁴ |
Julia Onken⁴ | Marcus Czabanka⁴ | Dongsheng Yuan^{1,5} | Sören Lukassen⁵ |
Philipp Karau⁵ | Naveed Ishaque⁵ | Elisabeth G. Hain⁶ | Frank Heppner⁶ |
Ahmed Idbaih⁷ | Nikolaus Behr¹ | Christoph Harms^{1,8} | David Capper^{6,9} |
Philipp Euskirchen^{1,9} 

Oslo Nanopore



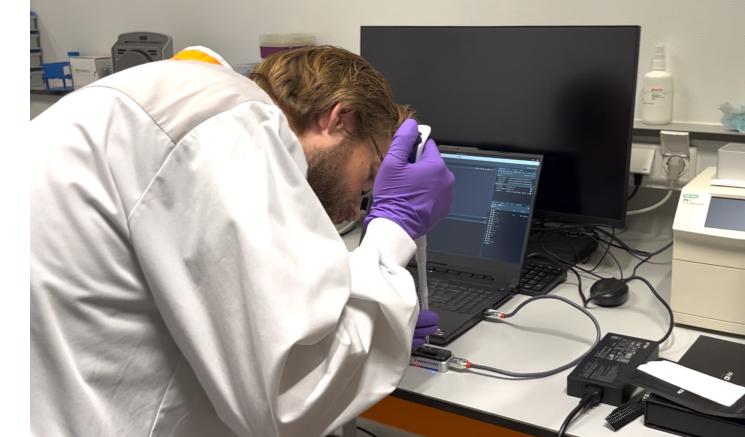
Real live run



Start
00:00



Finish DNA extraction
17:45



Start sequencing
26:00

```
2022-12-06 13:14:46.182 - root . INFO : Validating model: /home/sturgeon/nanocns/sturgeon_models/cv.256.128.adaptive/diagnostics_test/ensemble/model.zip
2022-12-06 13:14:46.183 - root . INFO : Successful model validation
2022-12-06 13:14:46.184 - root . INFO : Starting prediction session
2022-12-06 13:14:46.381 - root . INFO : Loading bed file: /home/sturgeon/nanocns/data/live_run_1/merged_probes_methyl_calls.bed
2022-12-06 13:14:46.382 - root . INFO : Total amount of measurable probes: 428643
2022-12-06 13:14:46.383 - root . INFO : Number of measured probes: 428643 (99.99%)
2022-12-06 13:14:46.384 - root . INFO : Number of measured methylated probes: 222 (5.37%)
2022-12-06 13:14:46.385 - root . INFO : Top 1: Sella - PITAD FSH LH - FSH LH (0.46%)
2022-12-06 13:14:46.386 - root . INFO : Top 2: Sella - PITAD FSH LH - FSH LH (0.46%)
2022-12-06 13:14:46.387 - root . INFO : Top 3: Pitadl - FSH LH - FSH LH (0.05%)
2022-12-06 13:14:46.388 - root . INFO : Saving results to: /home/sturgeon/nanocns/data/live_run_1/merged_probes_methyl_calls_model.csv
2022-12-06 13:14:46.389 - root . INFO : Plotting results to: /home/sturgeon/nanocns/data/live_run_1/merged_probes_methyl_calls_model.pdf
[1] "no new files, waiting 30 seconds"
```

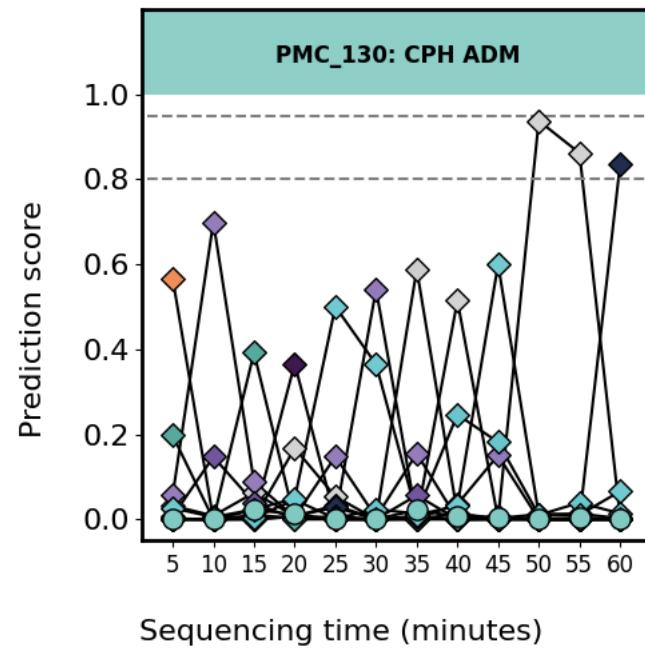
First prediction result
40:25



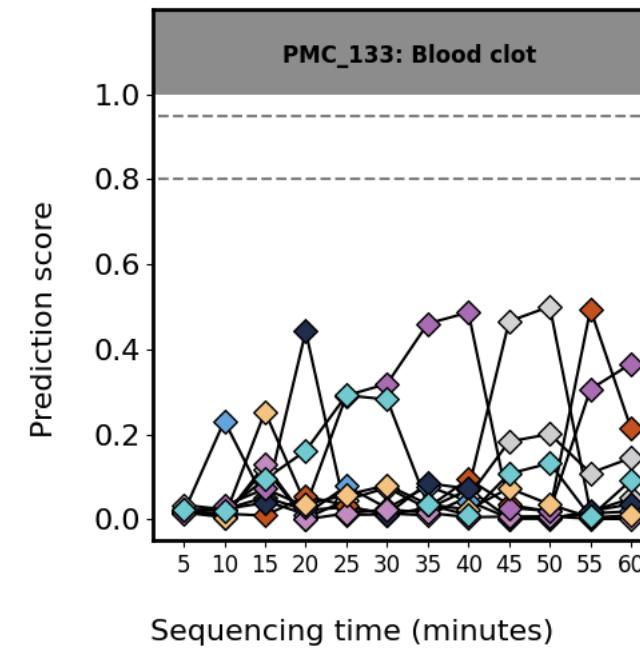
Doable in
one hour

Real live run

5 % sample purity

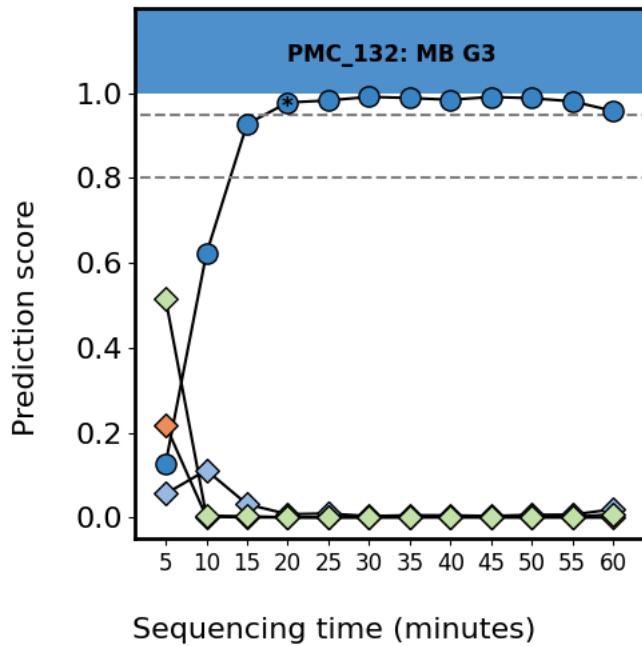


Blood clot

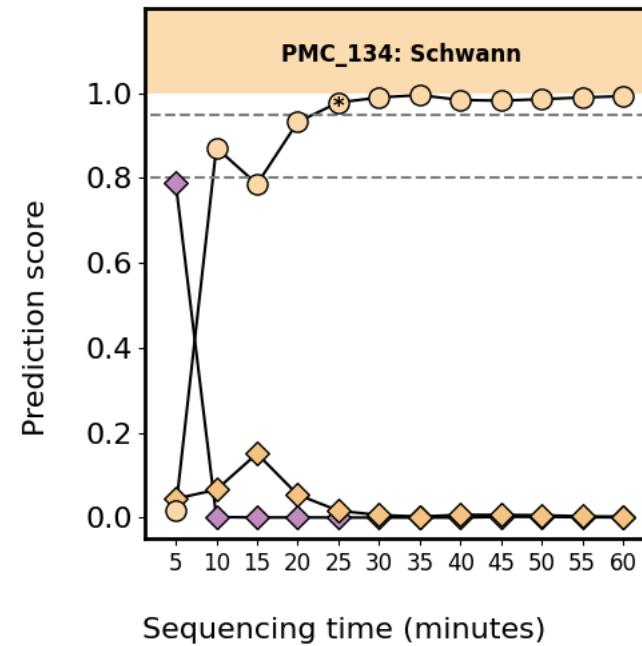


Real live run

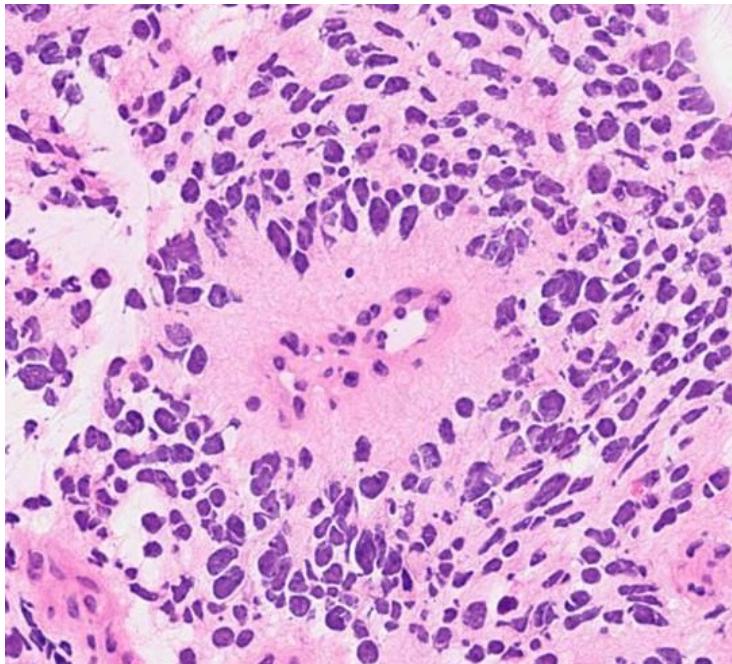
Medulloblastoma G3



Schwanomma



Real live run

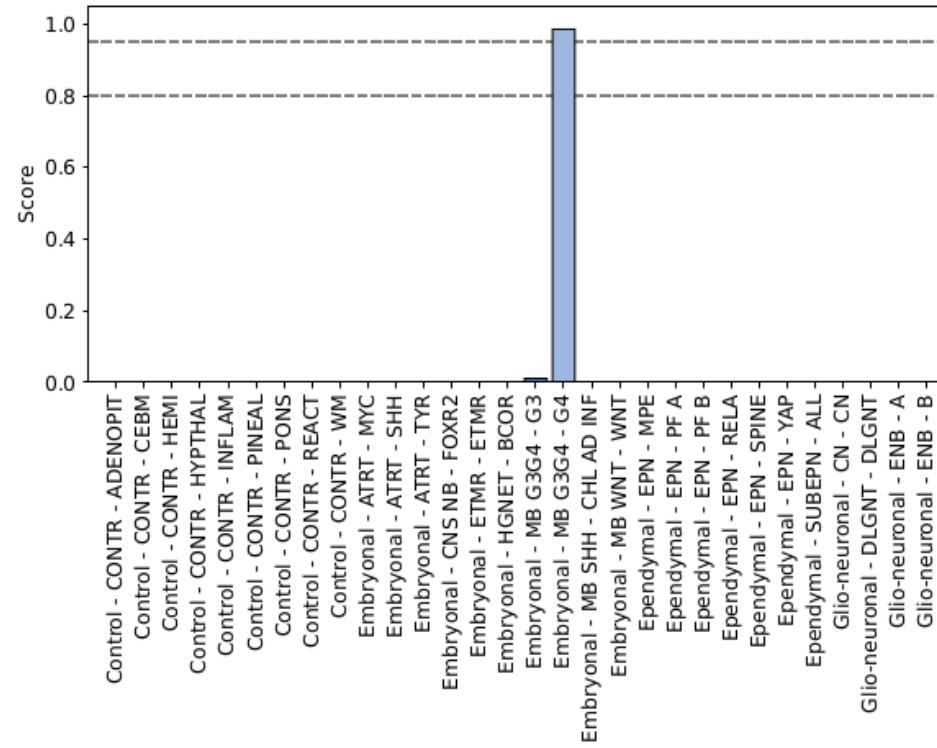


Histology: Ependymal Rosettes

Ependymoma: radical resection

One week later: methylation classifier says it's a Medulloblastoma G4

Sturgeon after 10 minutes of sequencing: **Medulloblastoma G4**



Conclusions

- We can train a neural network that can handle an arbitrary amount of random missing values and accurately predict brain cancer subtypes
- Data simulation allows us to
 - Massively increase sample size
 - Upsample to balance and also to improve on worse performing cases
 - Translate from microarray to Nanopore sequencing data
 - “Infinitely validate”
- All together is fast enough to be done and influence surgical strategy

Acknowledgements

de Ridder group

Jeroen de Ridder

Carlo Vermeulen

Adrien Melquiond

Dieter Stoker

Emmy Wesdorp

Inez den Hond

Joanna Wolthuis

Li-Ting Chen

Luca Santuari

Lucía Barbadilla

Nicolle Besselink

Roy Straver

Princess Maxima Center

Bastiaan Tops

Lennart Kester

Eric Strengman

Peter Wesseling

Eelco Hoving

Mariëtte Kranendonk

Kirsten van Baarsen

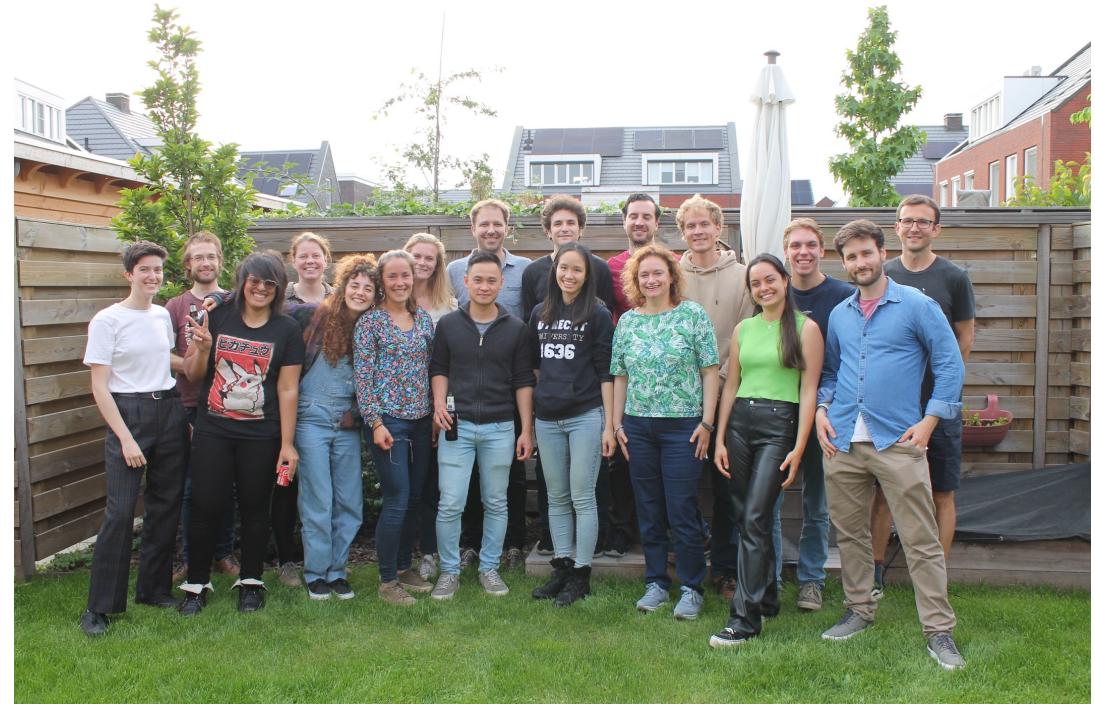
Jasper van der Lught

Universitätsklinikum

Münster

Christian Thomas

Anne Albers



UMC Utrecht



Oncode
Institute

Preprint available

de Ridder group

Jeroen de Ridder

Carlo Vermeulen

Adrien Melquiond

Dieter Stoker

Emmy Wessendorp

Inez den Hond

Joanna Wolthuis

Li-Ting Chen

Luca Santuari

Lucía Barbadilla

Nicolle Besselink

Roy Straver

Princess Maxima Center

Bastiaan Tops

Lennart Kester

Eric Strengman

Peter Wesseling

Eelco Hoving

Mariëtte Kranendonk

Kirsten van Baarsen

Jasper van der Lugt

Universitätsklinikum

Münster

Christian Thomas

Anne Albers



UMC Utrecht



Oncode
Institute