

RNA methylation calling with nanopore sequencing

Sebastian Krautwurst
RNA Bioinformatics
FSU Jena

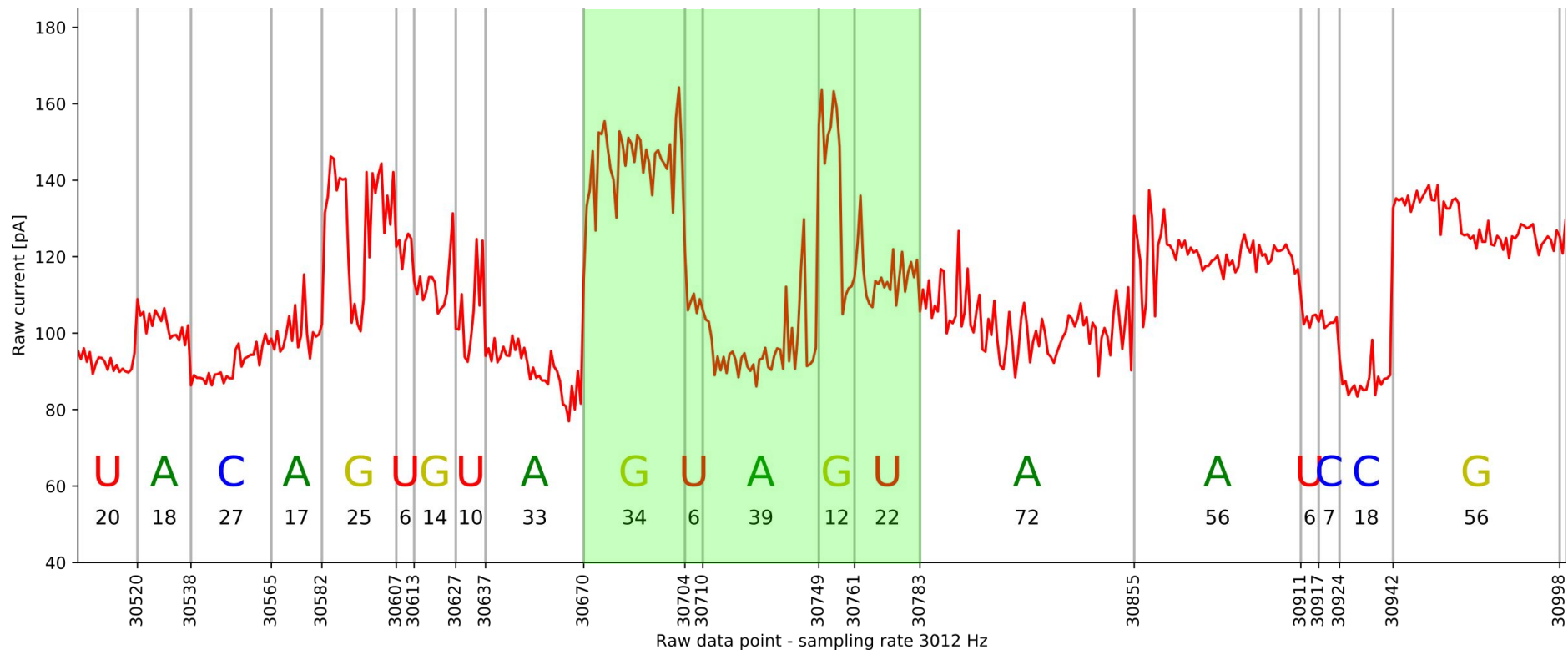
2020-02-11
35th TBI Winterseminar in Bled

Direct RNA sequencing with nanopores

- Protocol for MinION
- Directly sequences RNA strands
- No fragmentation
- No amplification
- Capture full transcripts
- Modifications are retained

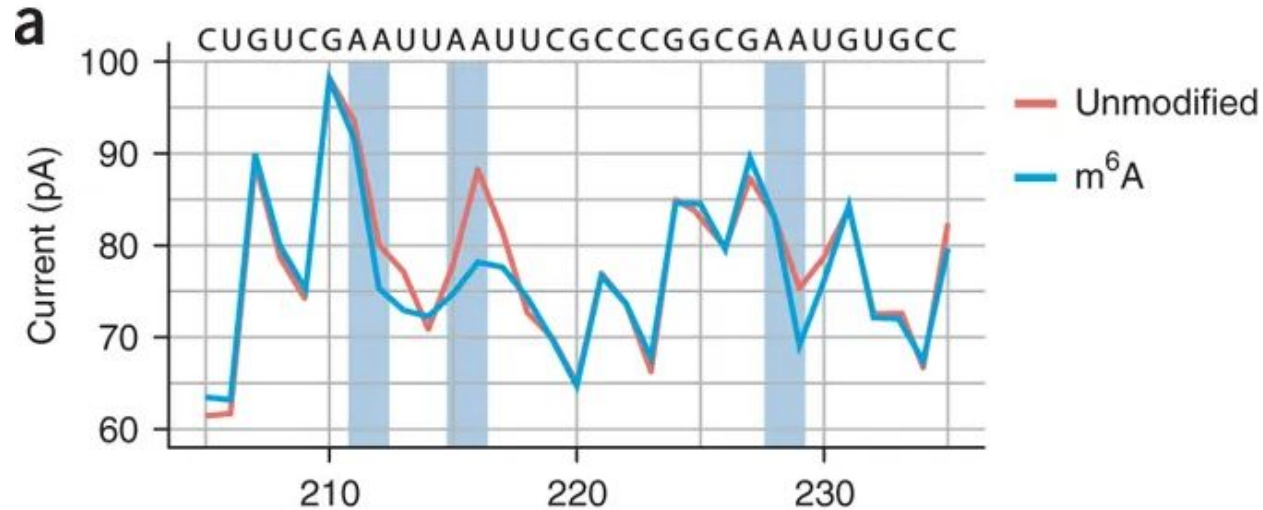


Basecalling: squiggle to basecalls

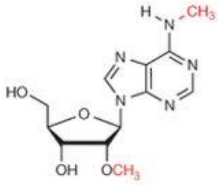
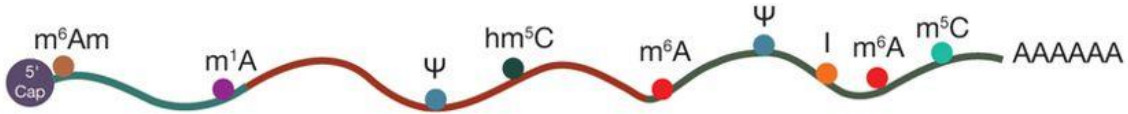


Modification calling

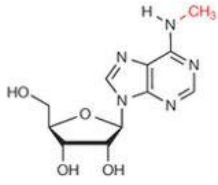
- Modifications change the signal
- DNA:



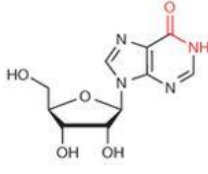
RNA modifications



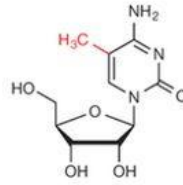
*N*⁶,2'-O-dimethyladenosine (*m*⁶Am)



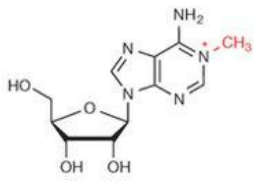
*N*⁶-methyladenosine (*m*⁶A)



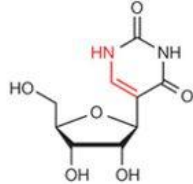
Inosine (I)



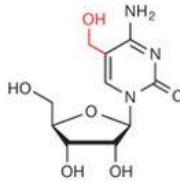
5-methylcytidine (*m*⁵C)



*N*¹-methyladenosine (*m*¹A)



Pseudouridine (Ψ)



5-hydroxymethylcytidine (*hm*⁵C)

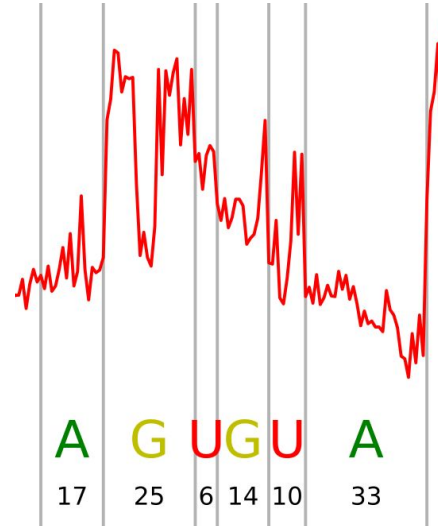
Guppy basecaller + Taiyaki training tool

Guppy:

- For DNA: 5mC, m6A calling is available
- Only trained for certain contexts: CpG, RRACH motif
- Basecall output + Modification output

Taiyaki:

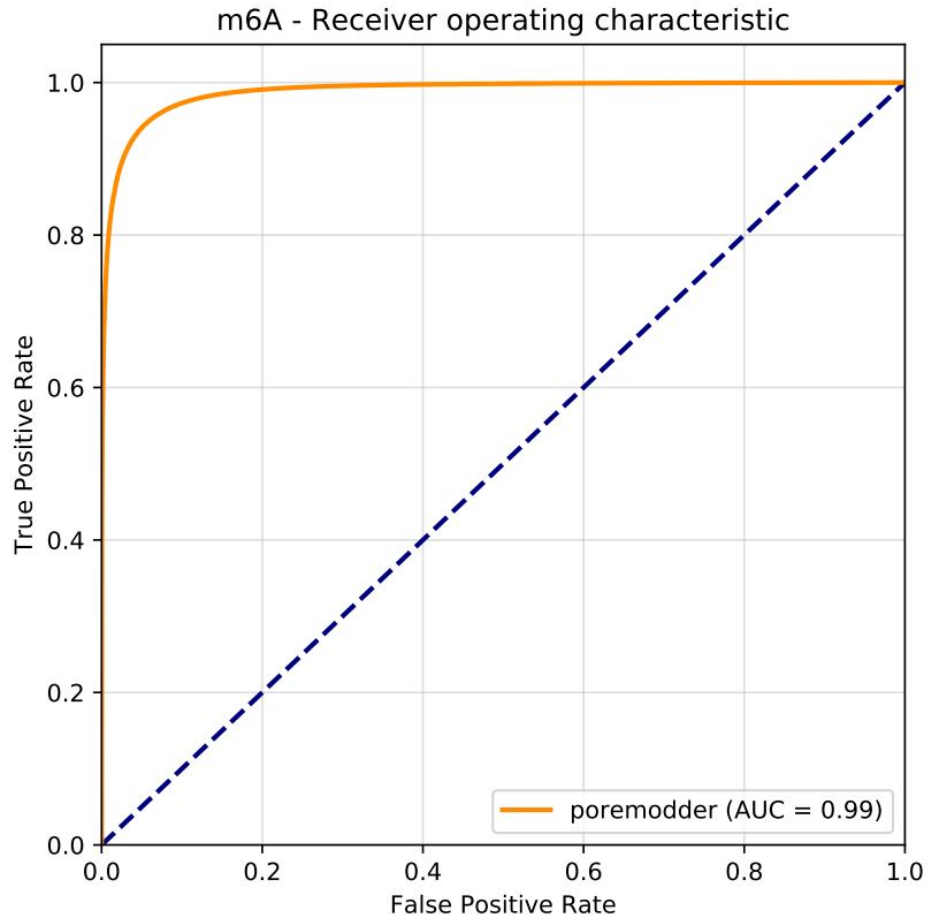
- Tool used for training Guppy
- Can train for modified bases
- Has “re-squiggle” functionality to prepare data for training

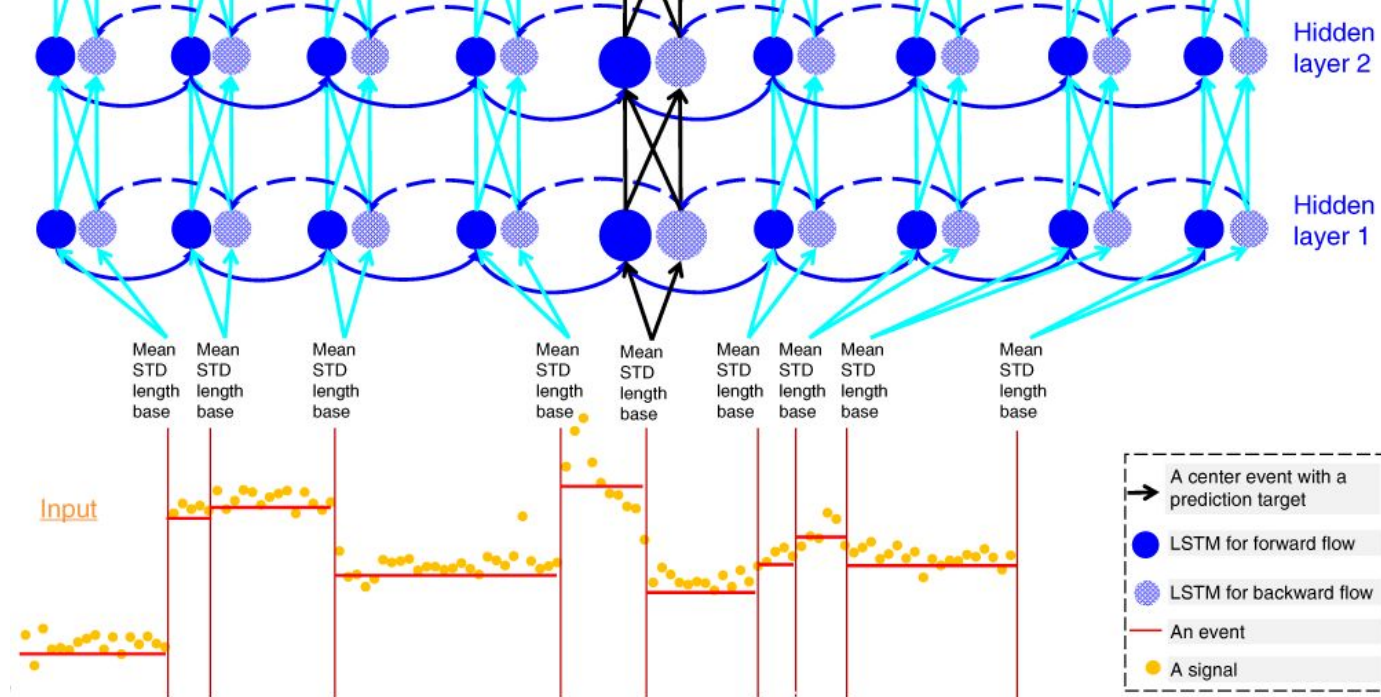


Training data for RNA?

- No good ground truth datasets
- In vitro transcripts:
 - With m6A - fully methylated
 - With A - fully canonical
- Model trained with Taiyaki

- Bad overfitting on training data
- Cannot basecall anything else





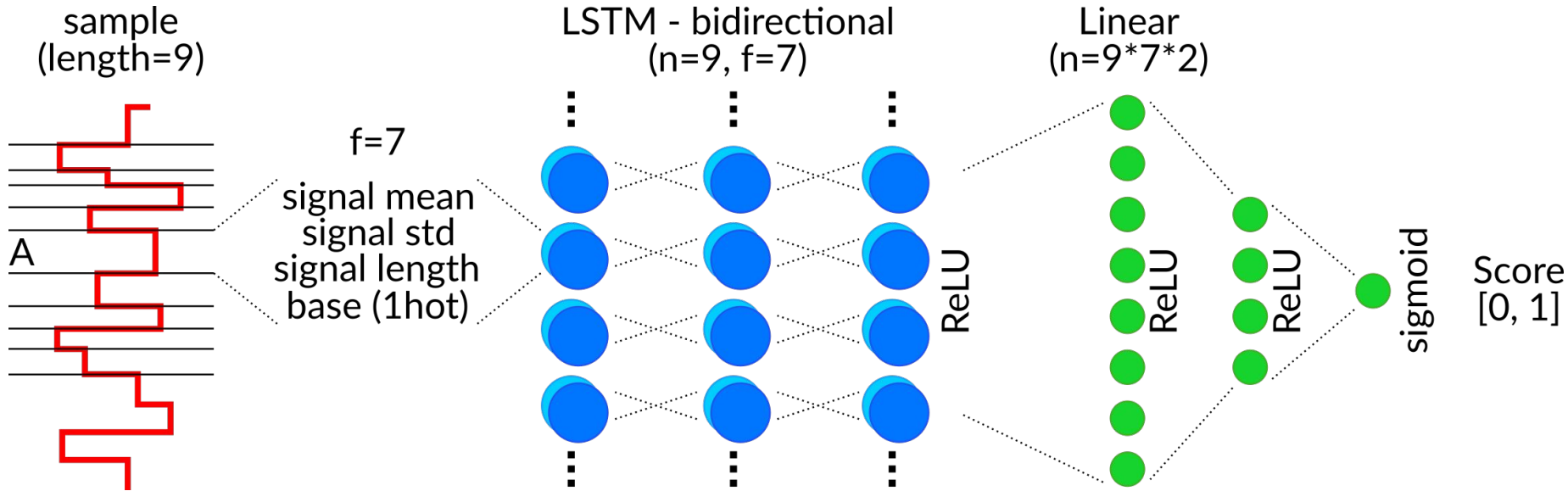
E. coli DNA data, per position:

- **5mC** 0.99 average precision
- **m6A** 0.9 average precision

Datasets - In vitro transcripts

m6A data	# reads canonical	# reads m6A	# samples (7mers)
Epinano replicate 1	46,628	8,739	15,990,176
Epinano replicate 2	603,275	94,046	TBD
Our own m6A data	in progress	26,621	TBD

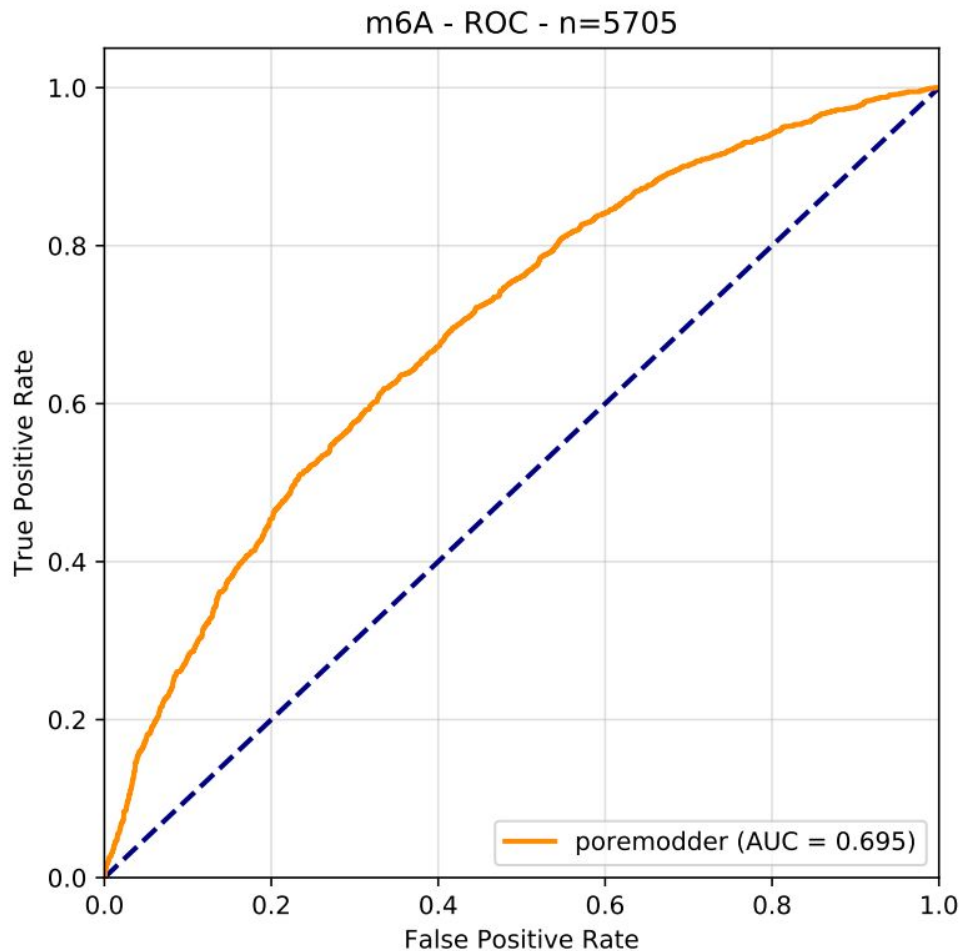
Prototype architecture using PyTorch



Results

- Trained on ~30k samples (9mers)
- Training/test split: 90/10
- Accuracy: 0.661

	Prec.	Recall	F1-score	#
Nomod	0.819	0.687	0.747	4157
m6A	0.413	0.592	0.487	1548



Results

```
[1, 3409500] loss: 0.244
      precision    recall  f1-score   support

   Nomod         0.668     0.674     0.671     778134
    m6A         0.638     0.631     0.634     707394

 accuracy                   0.653     1485528
 macro avg         0.653     0.652     0.653     1485528
weighted avg         0.653     0.653     0.653     1485528
```

Traceback (most recent call last):

```
File "train_lstm_model.py", line 307, in <module>
    means = data['means'].T.unsqueeze(2)
```

```
File "/home/ya86gul/miniconda3/envs/poremodder/lib/python3.7/site-packages/
    _error_if_any_worker_fails()
```

```
RuntimeError: DataLoader worker (pid 24510) is killed by signal: Killed.
```

Problems and Further Ideas

- Slow/bad data loading
- Data insufficient, e.g. k-mers with multiple As

- Data preprocessing?
- Data filtering?
- Feature selection?
- Network architecture?
- Activation functions?

Conclusion and Outlook

- Available data is problematic
- Task appears learnable

- Architecture prototyping
- Make more data!
- Punish those GPUs!





FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

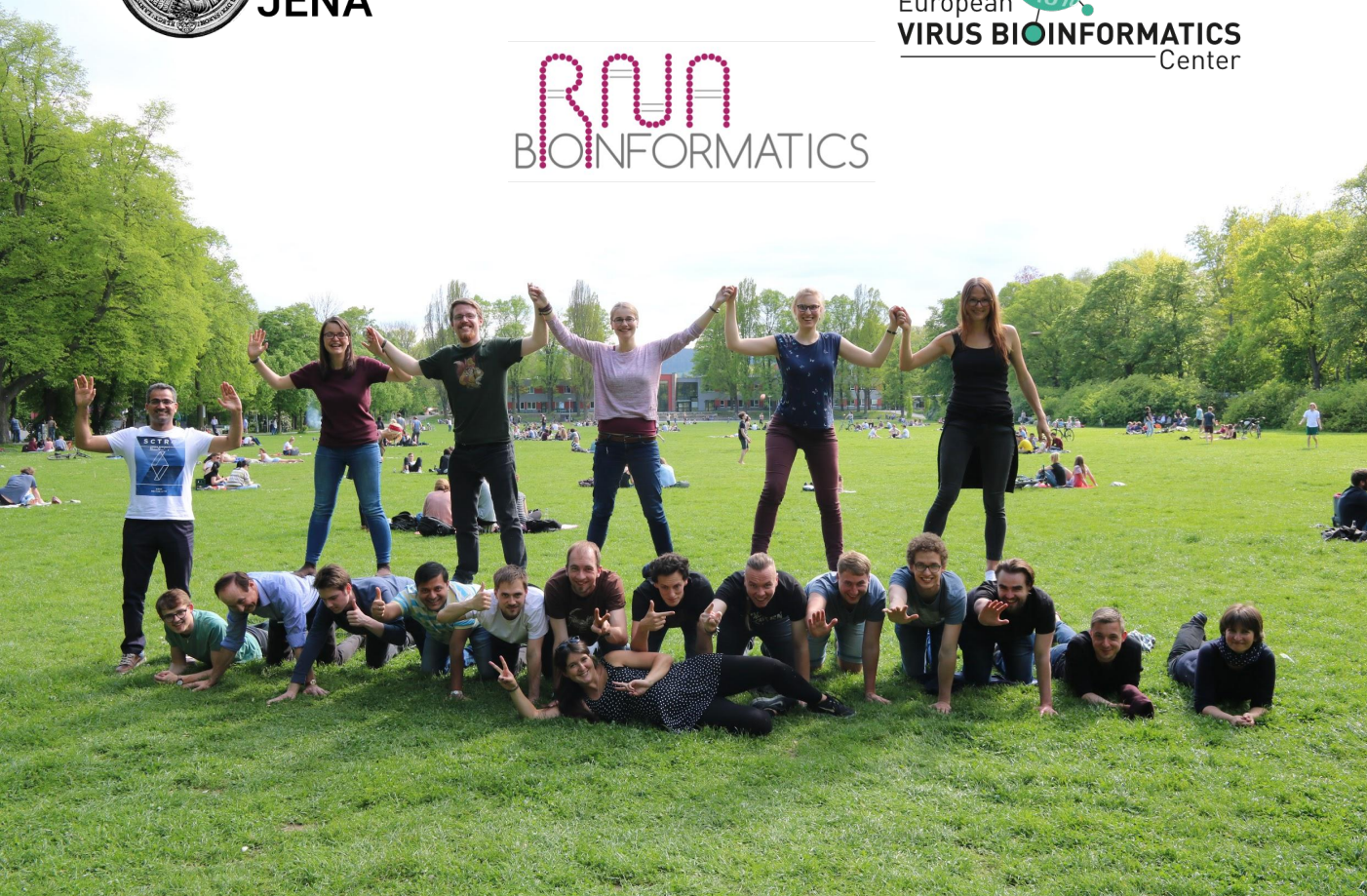


European
VIRUS BIOINFORMATICS
Center

RAUN
BIOINFORMATICS



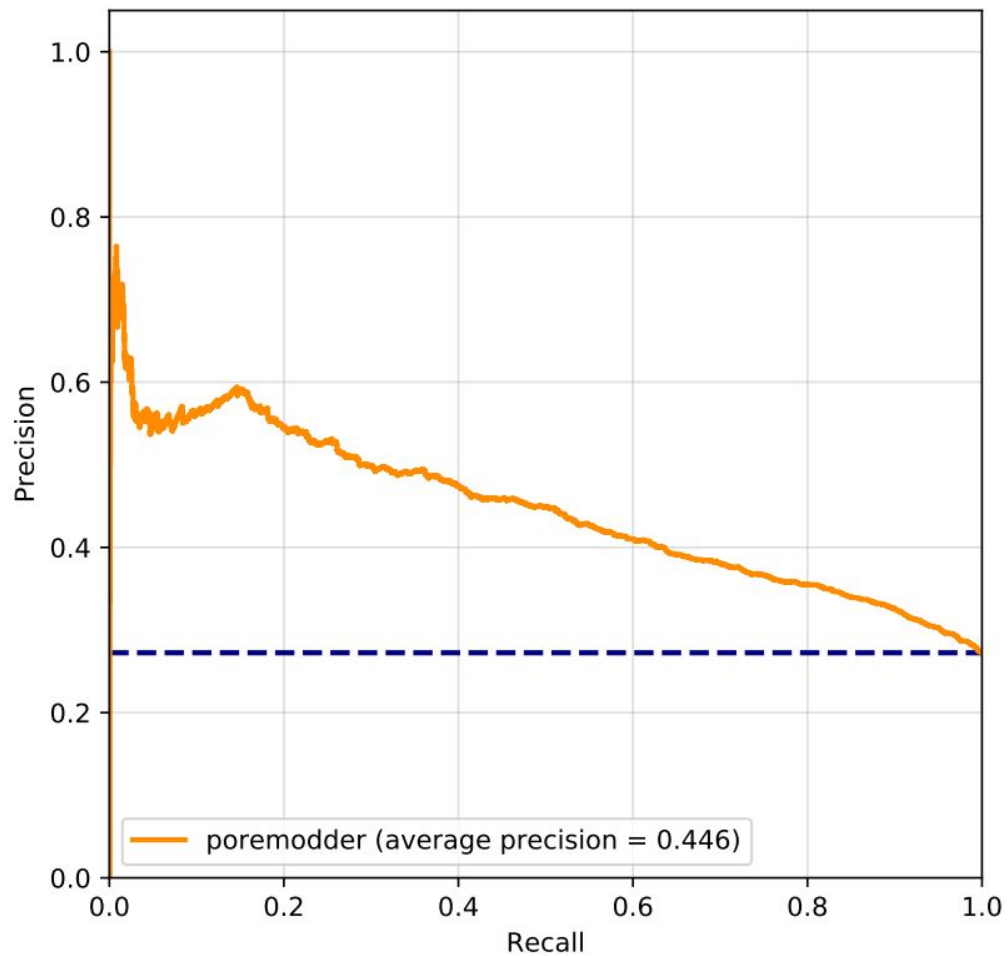
Celia Diezel Florian Mock
Manja Marz



Thank
you!

SK supported by
Oxford Nanopore
Technologies bursary

m6A - Precision-Recall - n=5705



sensitivity, recall, hit rate, or true positive rate (TPR)

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

specificity, selectivity or true negative rate (TNR)

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

precision or positive predictive value (PPV)

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

negative predictive value (NPV)

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$$

miss rate or false negative rate (FNR)

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

fall-out or false positive rate (FPR)

$$\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

false discovery rate (FDR)

$$\text{FDR} = \frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV}$$

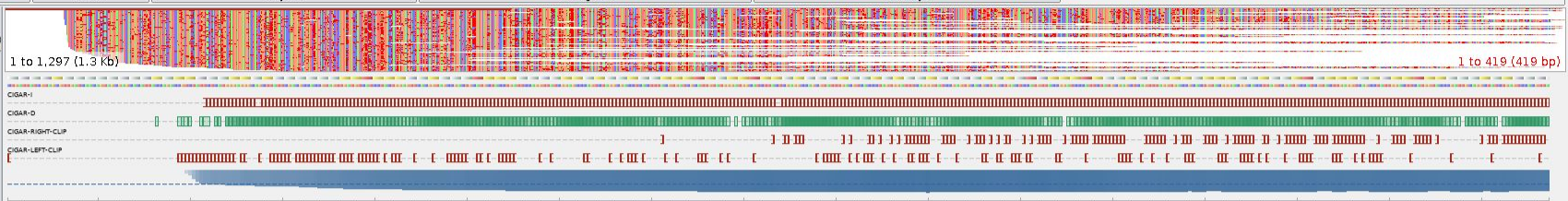
false omission rate (FOR)

$$\text{FOR} = \frac{\text{FN}}{\text{FN} + \text{TN}} = 1 - \text{NPV}$$

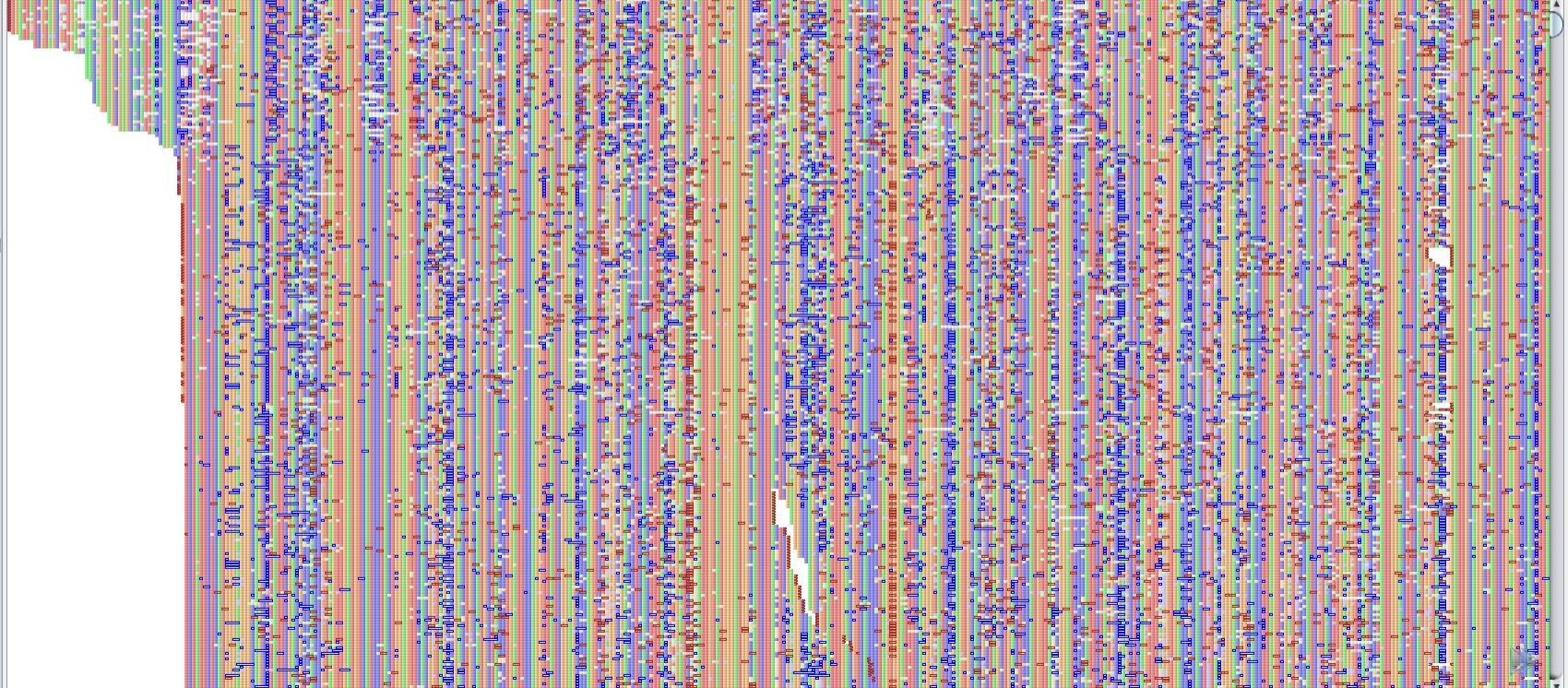
```
6 def network(insize=1, size=256, winlen=19, stride=2, alphabet_info=None):
7     return Serial([
8         Convolution(insize, size, winlen, stride=stride, fun=tanh),
9         Reverse(GruMod(size, size)),
10        GruMod(size, size),
11        Reverse(GruMod(size, size)),
12        GruMod(size, size),
13        Reverse(GruMod(size, size)),
14        GlobalNormFlipFlopCatMod(size, alphabet_info),
15    ])
```

Contigs: 1 26.78 k reads (more)

Co...	Le...	Re...	Fe...	Mis...
cu...	1,2...	26,...	4,8...	18,5...



1 U1 418 U418



Filter by:

<https://www.nvidia.com/content/dam/en-zz/Solutions/geforce/geforce-rtx-turing/2080/gallery/geforce-rtx-2080-gallery-c.jpg>

https://www.kxly.com/content/uploads/2019/12/generic-fire_1562686082475-jpg_38950375_ver1-0-1024x576.jpg