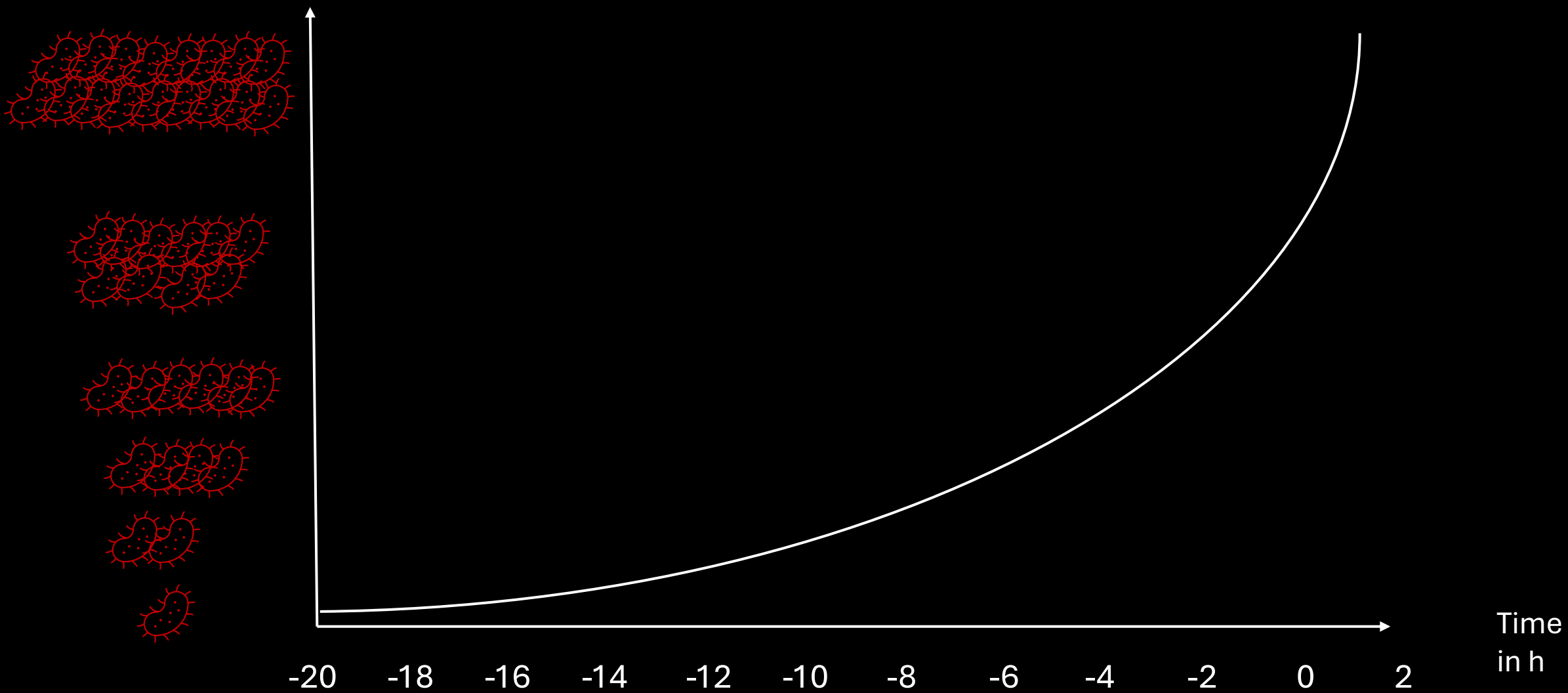
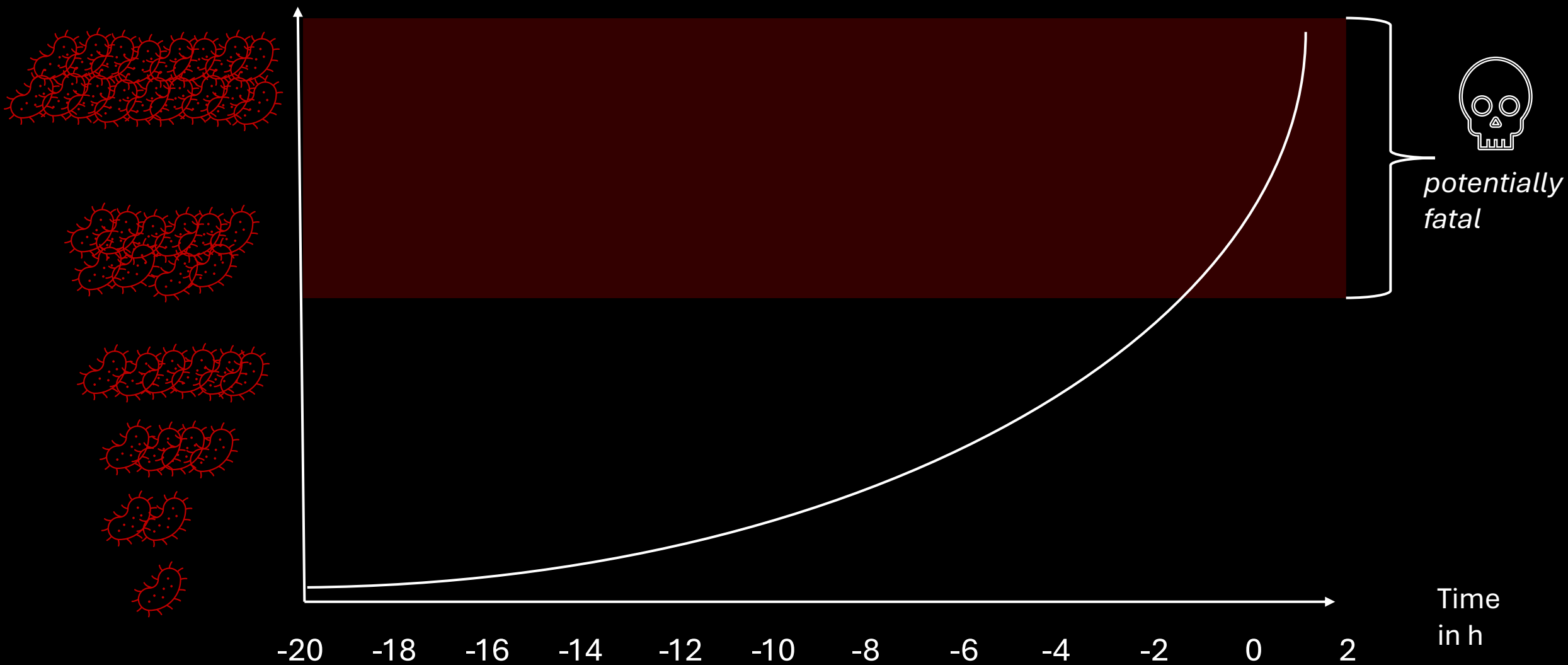
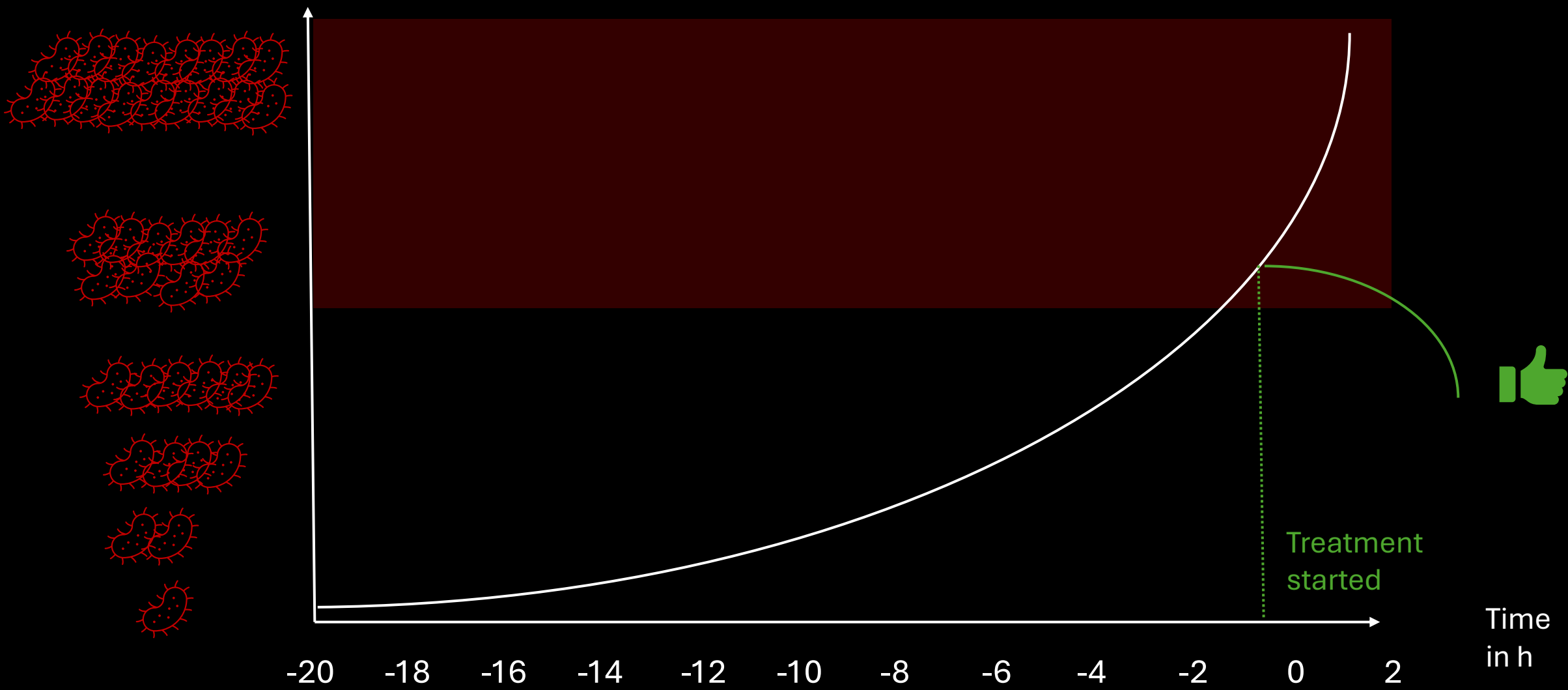


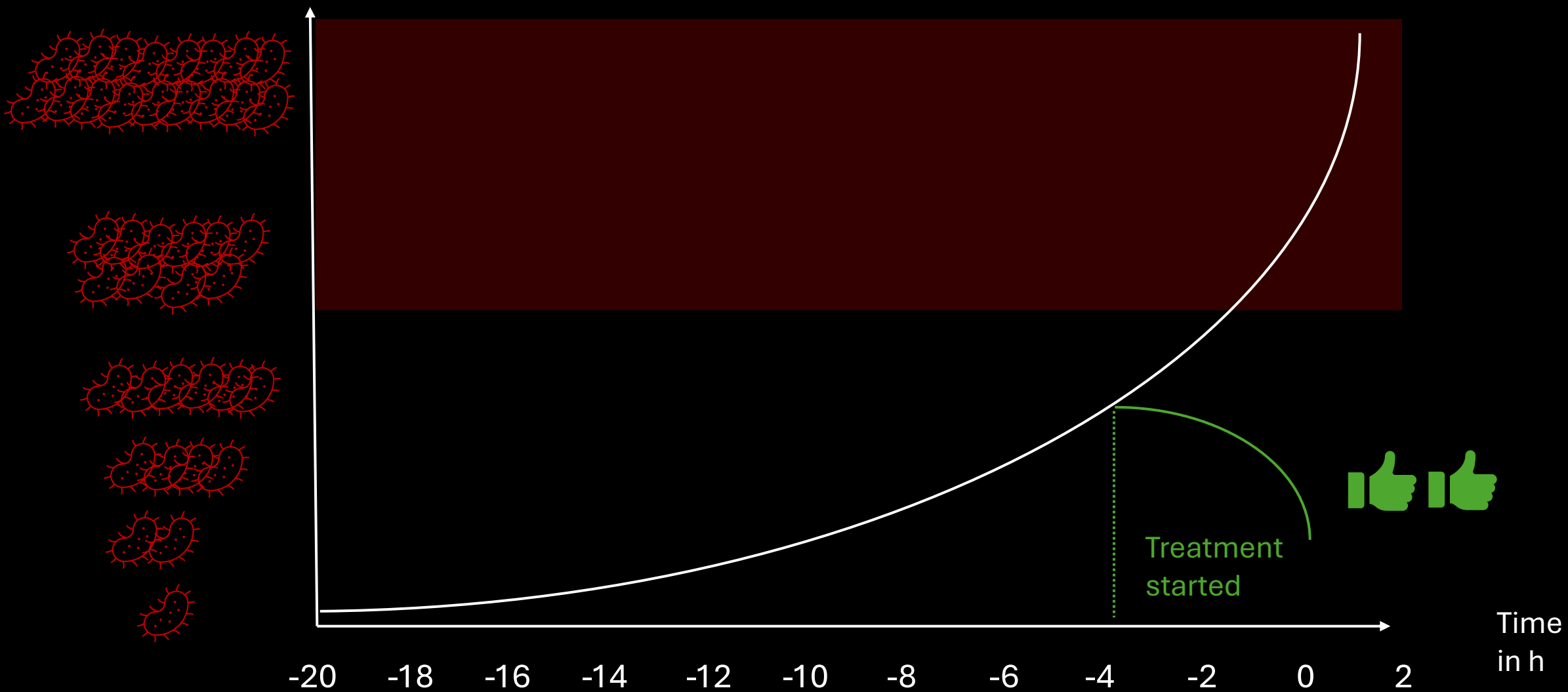


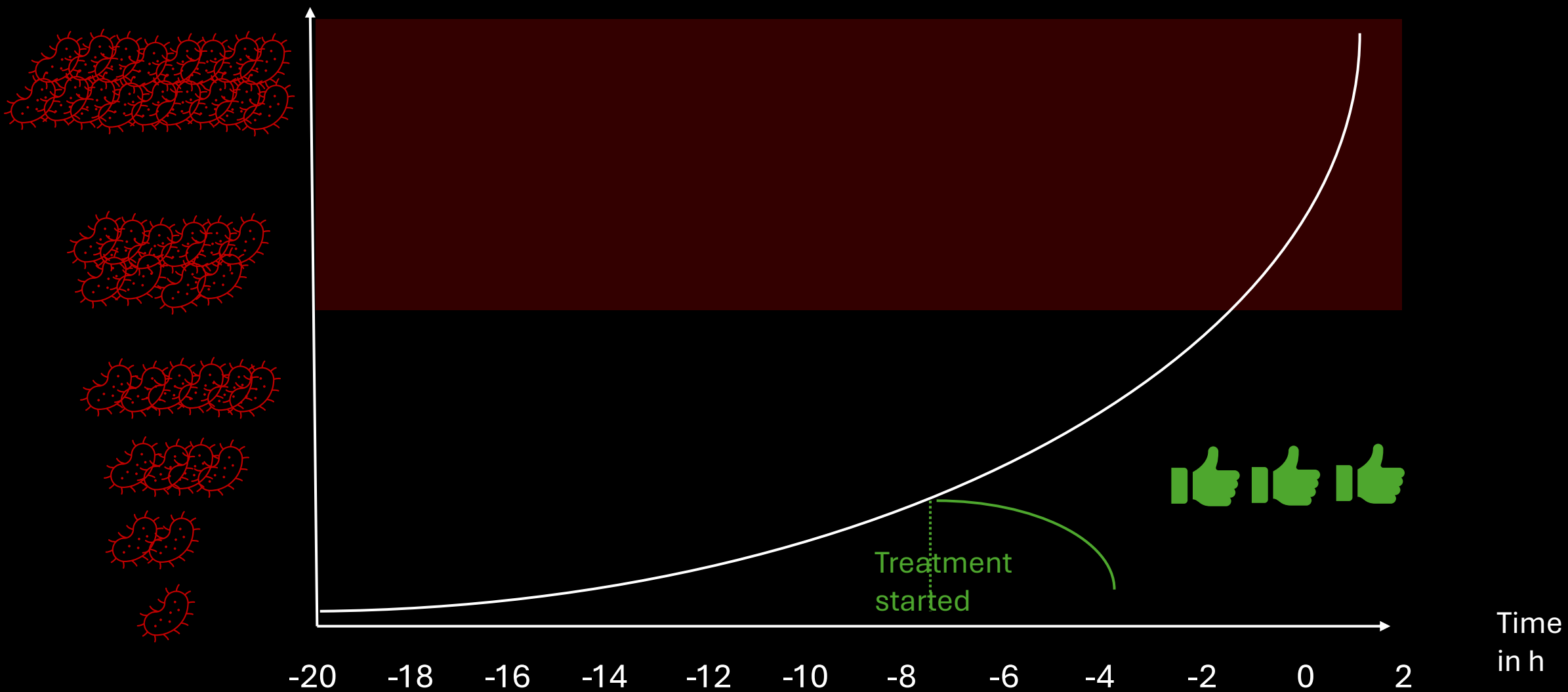
PREDICTING SEPSIS











“Detecting sepsis early and starting immediate treatment is often the difference between *life and death*.”

Problem Definition










Will the patient develop sepsis within the next 6 hours based on currently available data?

Yes or No

-> Binary classification task

Dataset

Early Prediction of Sepsis from Clinical Data: The PhysioNet/Computing in Cardiology Challenge 2019

Matthew Reyna , Chris Josef , Russell Jeter , Supreeth Shashikumar , Benjamin Moody , M. Brandon Westover , Ashish Sharma , Shamim Nemati , Gari D. Clifford 

Published: Aug. 5, 2019. Version: 1.0.0

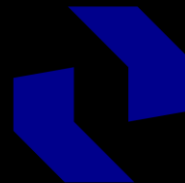


UNIVERSITY OF
OXFORD



UNIVERSITY OF
SOUTHERN CALIFORNIA

PHILIPS
Healthcare



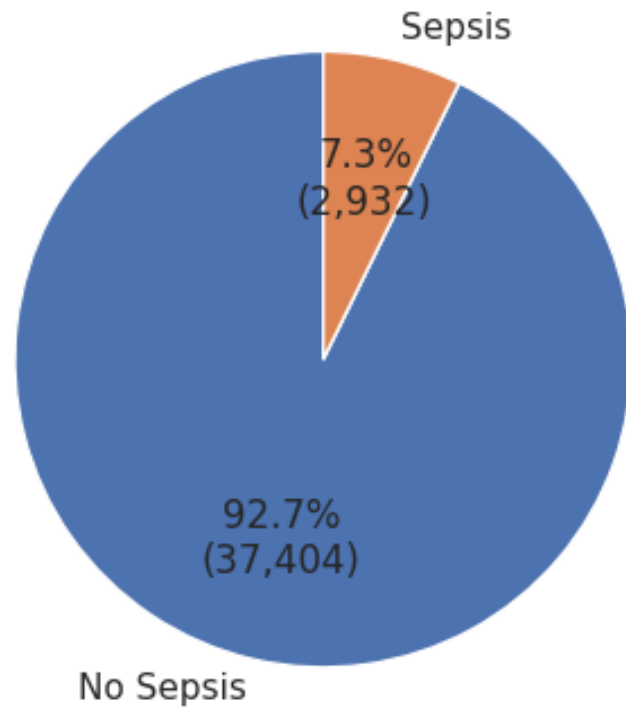
Technische
Universität
Dresden

and 55 more entries....

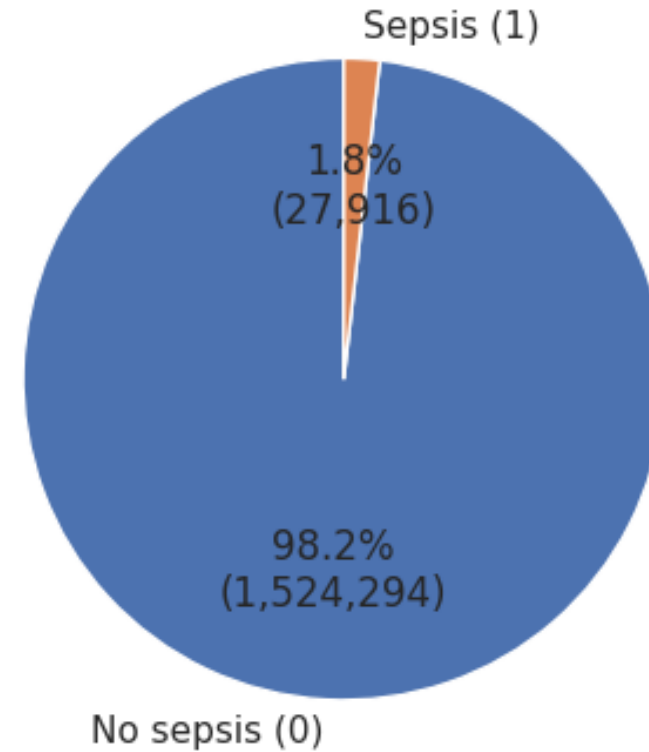
Patient 000001 as an example:

Dataset

Patient - wise Sepsis Distribution (N = 40,336)



Datapoint-wise Sepsis Labels



-> 1st challenge: Extreme Imbalance

Dataset

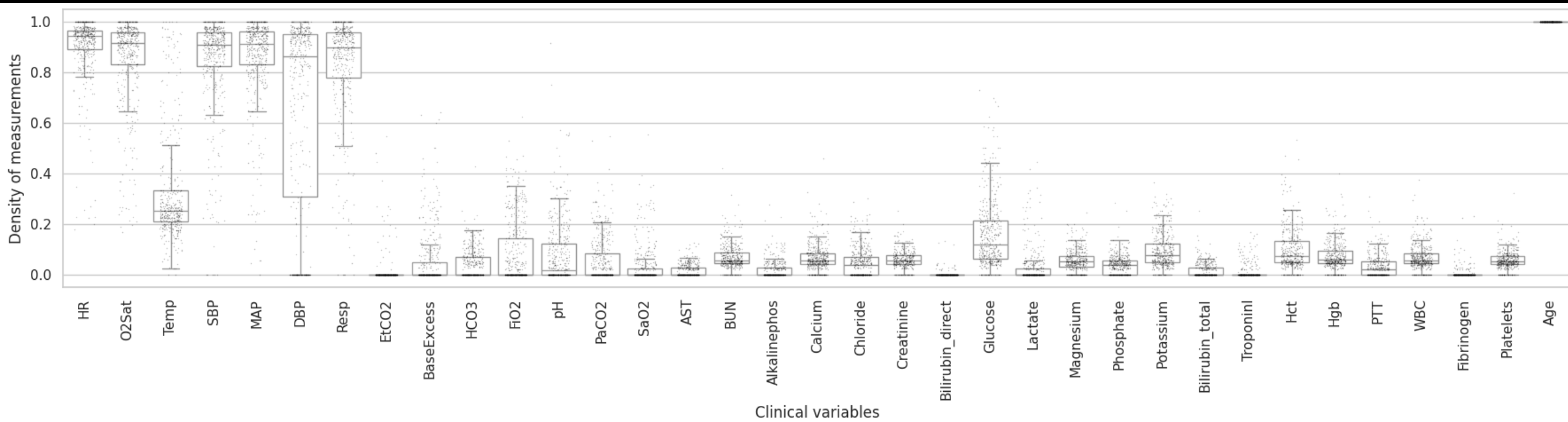
vital signs

Laboratory values

HR
O2Sat
Temp
SBP
MAP
DBP
Resp
EtCO2
BaseExcess
HCO3
FiO2
pH
PaCO2
SaO2
AST
BUN
Alkalinephos
Calcium
Chloride
Creatinine
Bilirubin_direct
Glucose
Lactate
Magnesium
Phosphate
Potassium
Bilirubin_total
TroponinI
Hct
Hgb
PTT
WBC
Fibrinogen
Platelets
Age

Clinical variables

Dataset



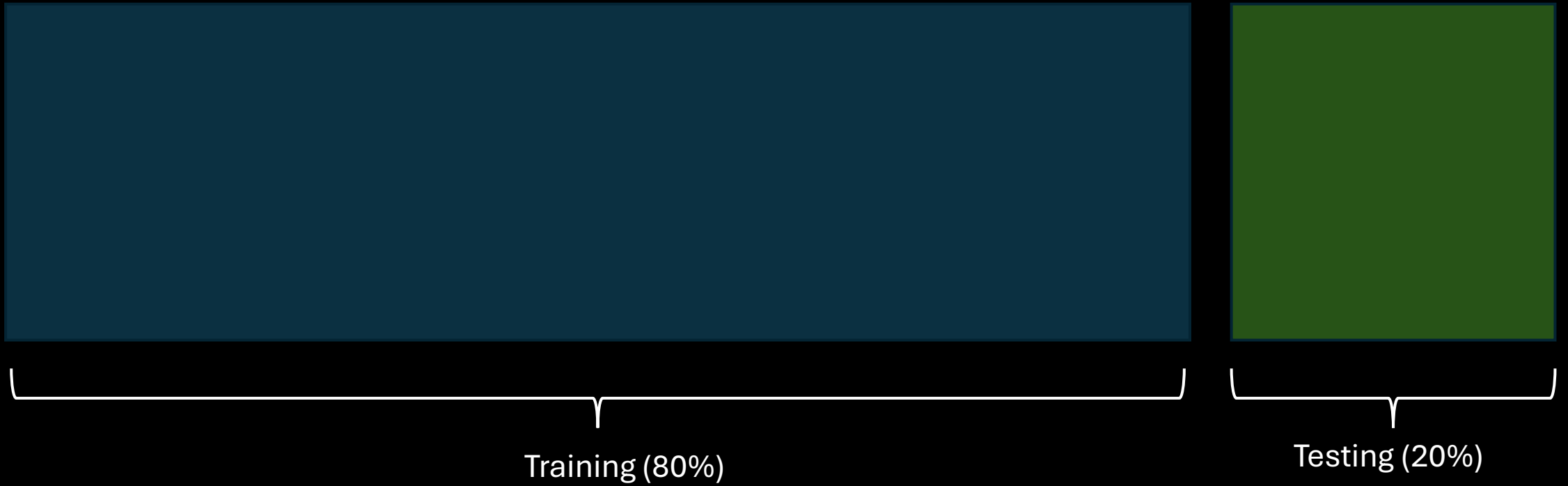
-> 2nd challenge = Many missing values!

Dataset

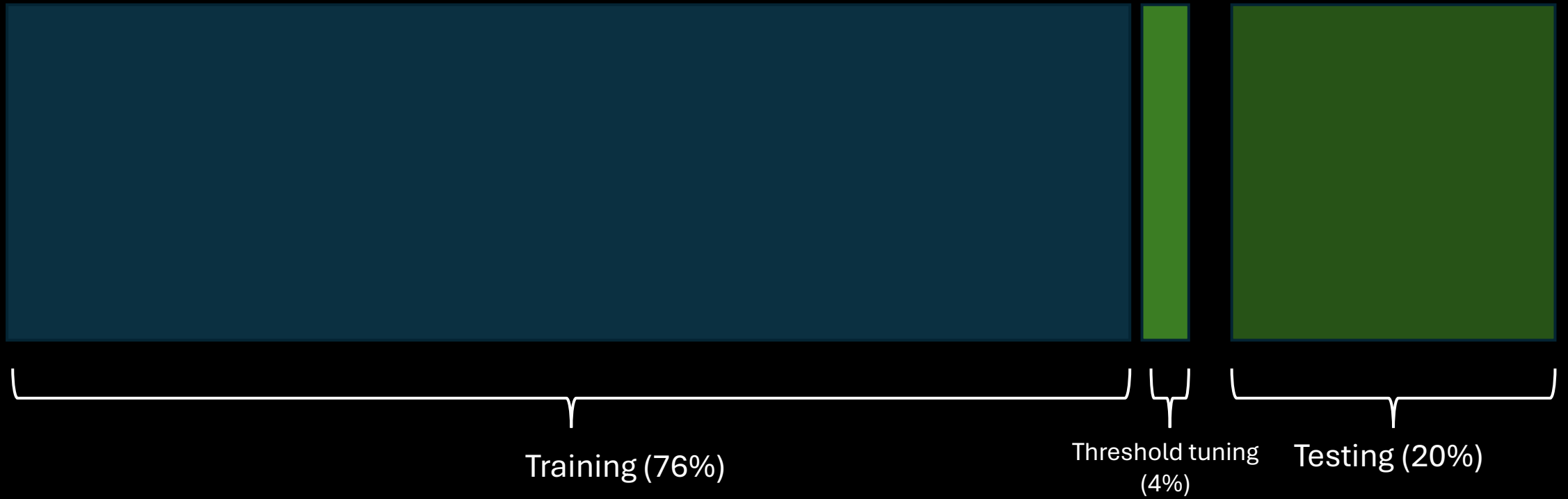
Full dataset: 40,336 patients



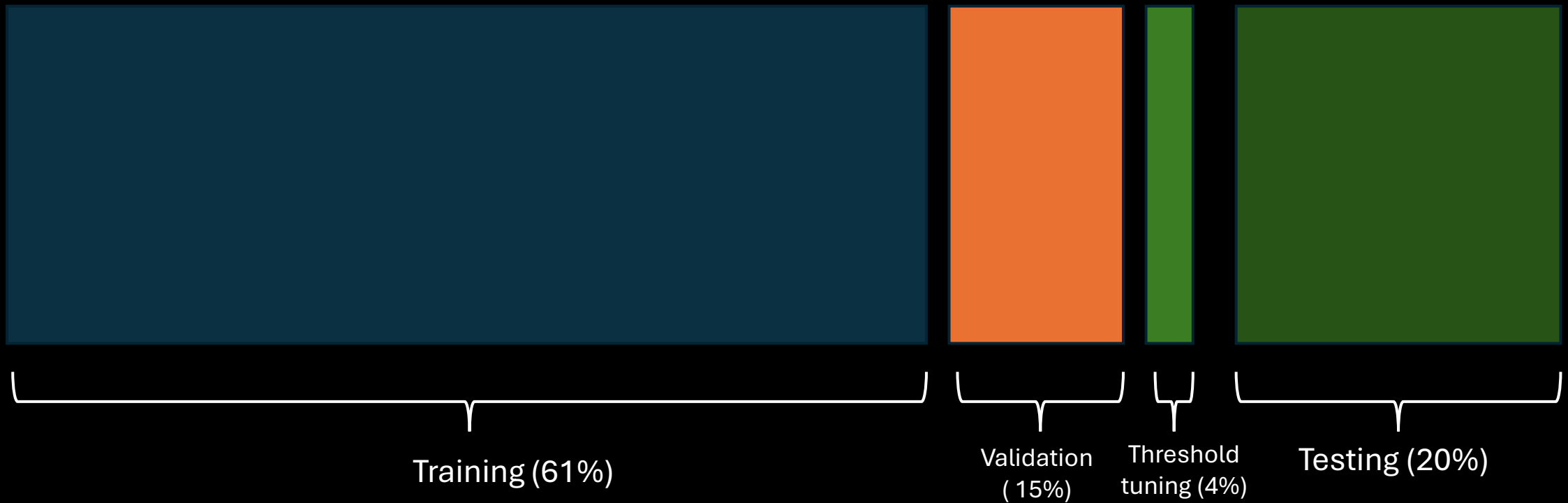
Dataset



Dataset



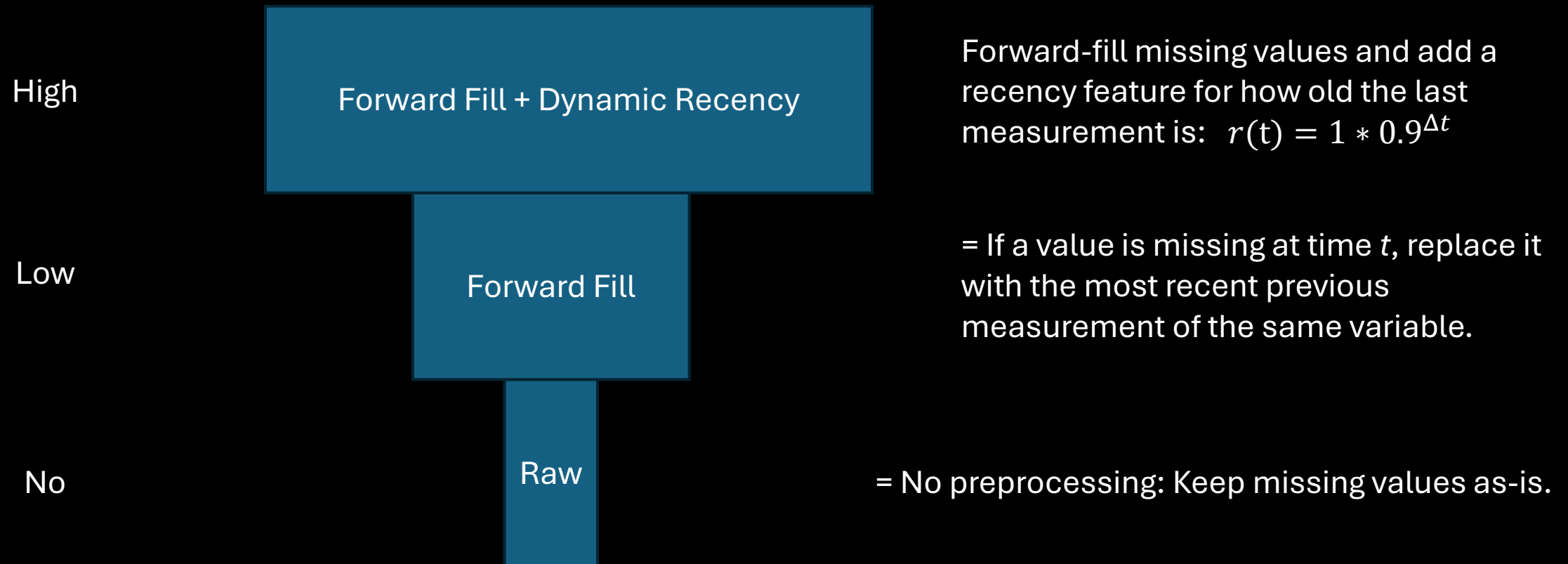
Dataset



-> All splits were done using a **stratified** split based on whether a patient ever gets septic

Handling Missing Data

Three preprocessing strategies were attempted:

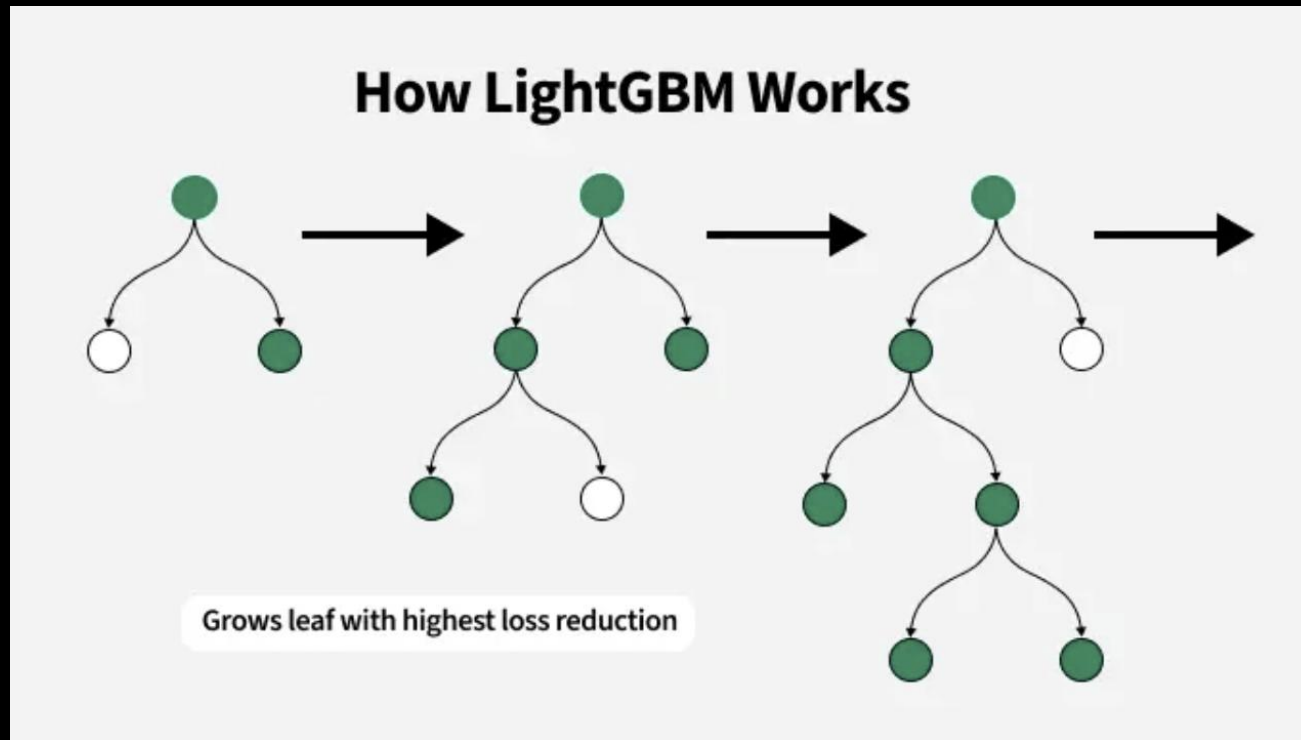


Baseline Model

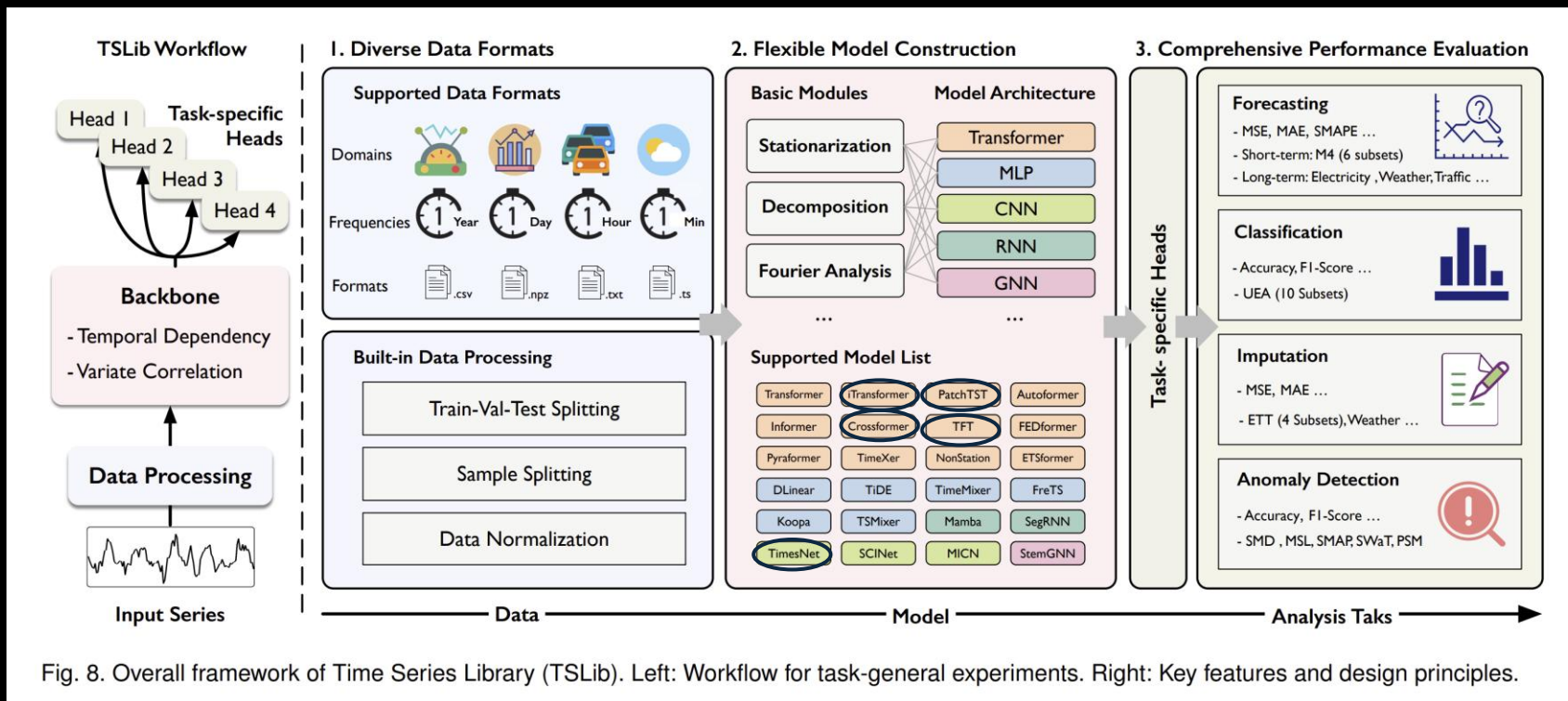
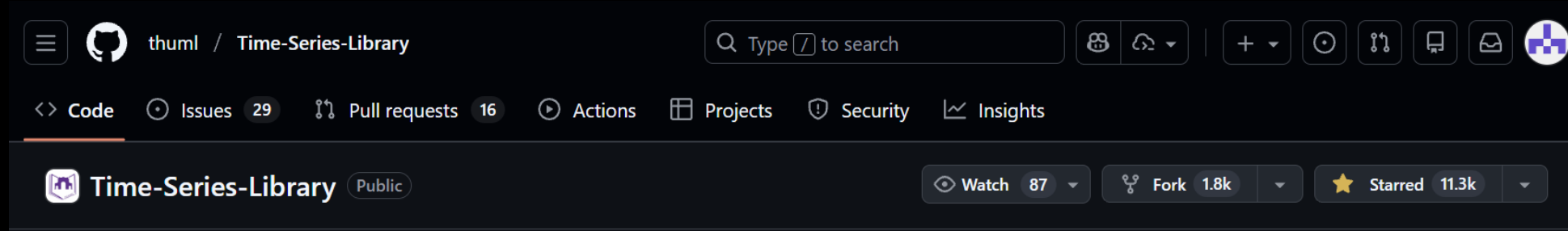
LightGBM

(Light Gradient-Boosting Machine)

- ➔ Efficient and scalable
- ➔ Robust and interpretable
- ➔ Proven to perform well on tabular data (Most popular choice in PhysioNet 2019 Competition)



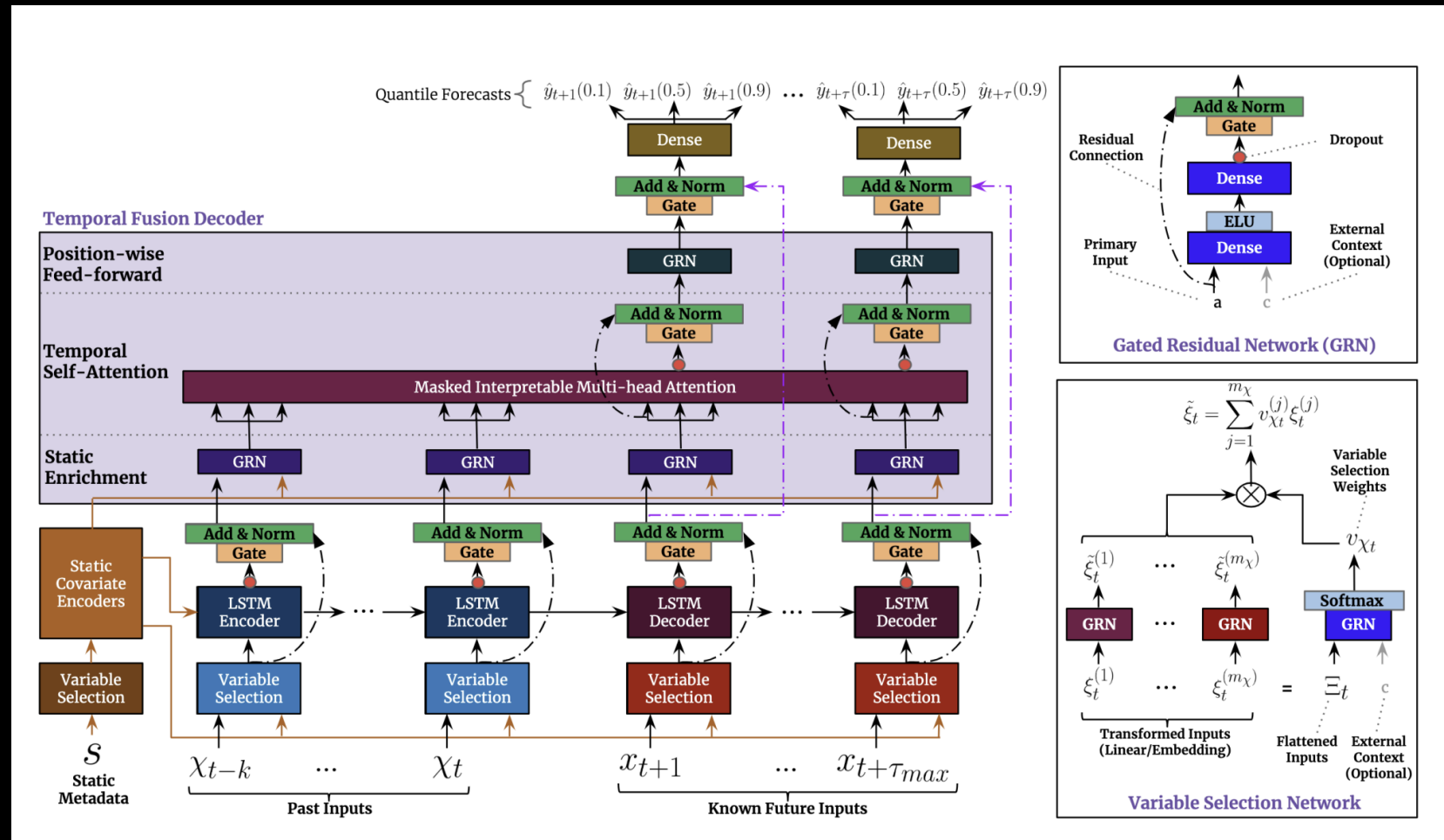
Deep Learning Models: Framework



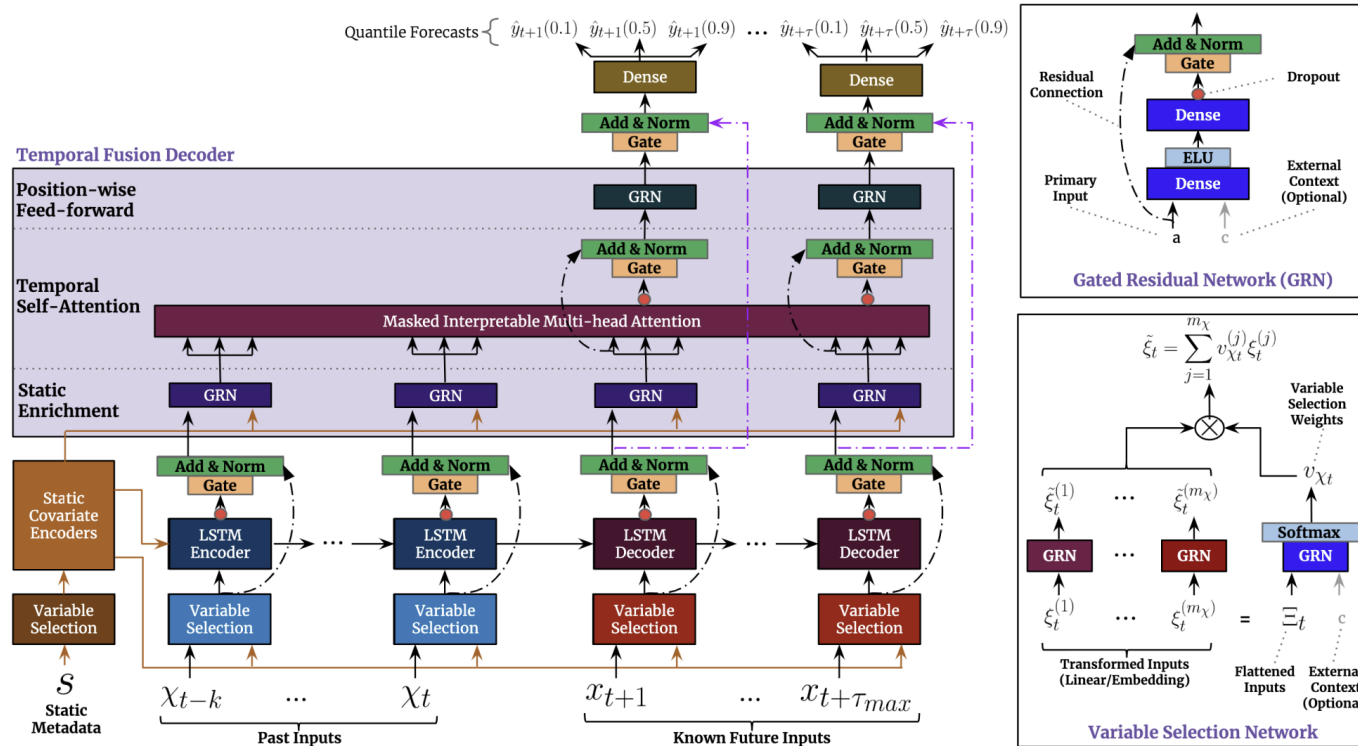
“A Library for Advanced Deep Time Series Models for General Time Series Analysis.”

- ➔ Fully open source 🍏
- ➔ Modular approach 🍏

Deep Learning Models: Temporal Fusion Transformers



Deep Learning Models: Temporal Fusion Transformers

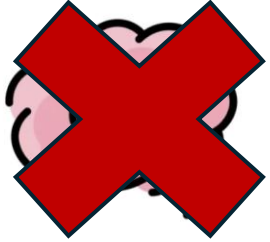


- ➔ Designed for complex time series with heterogeneous inputs
- ➔ Captures both short- and long-term temporal dependencies
- ➔ High performance with built-in interpretability

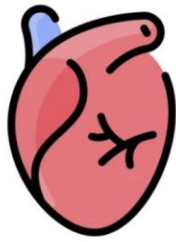
Clinical Screening Method: partial qSOFA

qSOFA-Score

Quick Sepsis-Related Organ Failure Assessment



GCS < 15 Pkt



RR sys \leq 100 mmHg



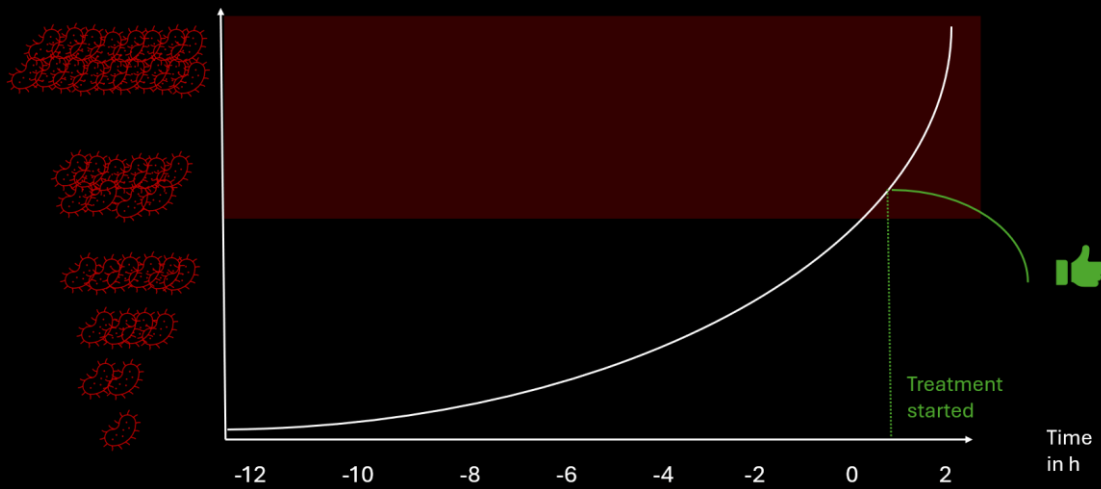
AF \geq 22/min

Problem: GCS not in dataset,
therefore only *partial* qSOFA

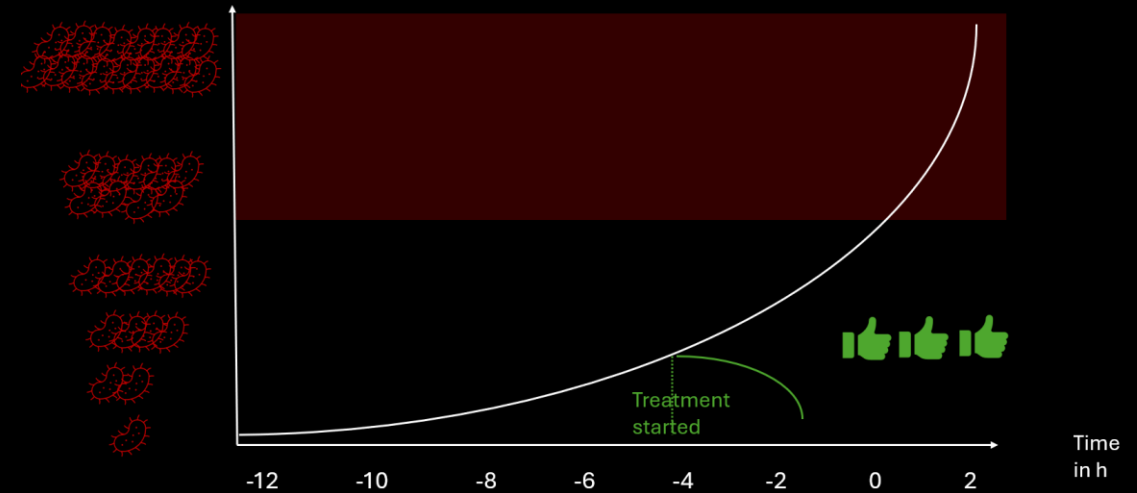
qSOFA is the most commonly used screening tool for sepsis

-> But: only based on clinical (i.e. non-laboratory) information.

Evaluation

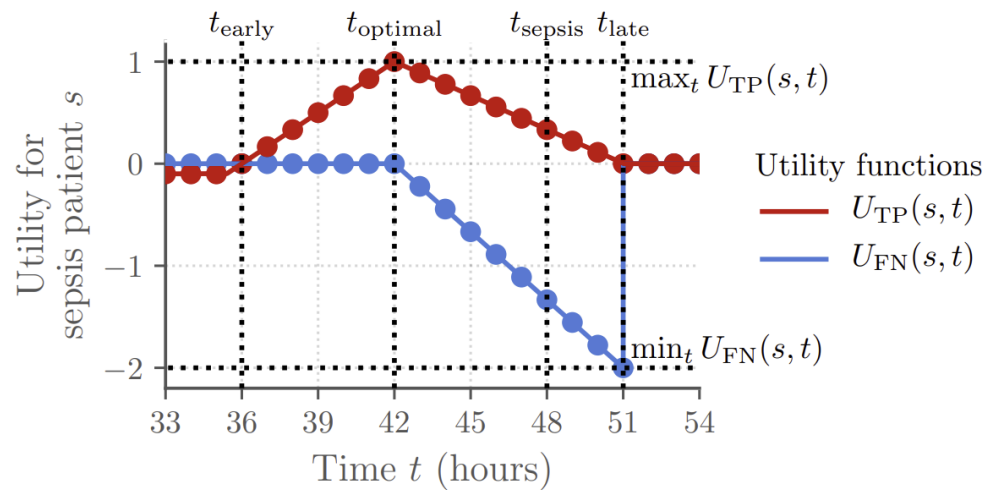


VS

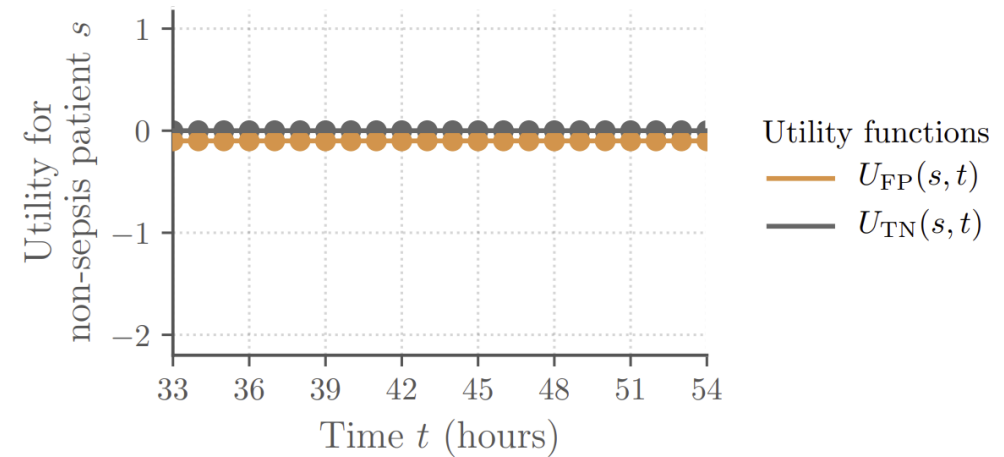


➔ How do we *reward* early sepsis predictions more than late ones and how do we *penalize* missed ones adequately?

Evaluation: Utility Score



(a) Clinical utility for septic patients.



(b) Clinical utility for non-septic patients.

Evaluation

Model(s)
create
hourly
predictions

Patient 0000007 as Example

1	PredictedProbability
2	0.5340326920919986
3	0.7573166166875013
4	0.7970109585800176
5	0.7217623506906001
6	0.680362839730447
7	0.7370480555111093
8	0.7654662167703292
9	0.7606485477998423
10	0.7597913835227376
11	0.6931276372372849
12	0.7263450231403222
13	0.7163432308875632
14	0.7129223634757045
15	0.699253567759149
16	0.7015488057939573
17	0.7120288374607344
18	0.7240207015205272



Standardized
Threshold sweep
was conducted
to find the best
threshold

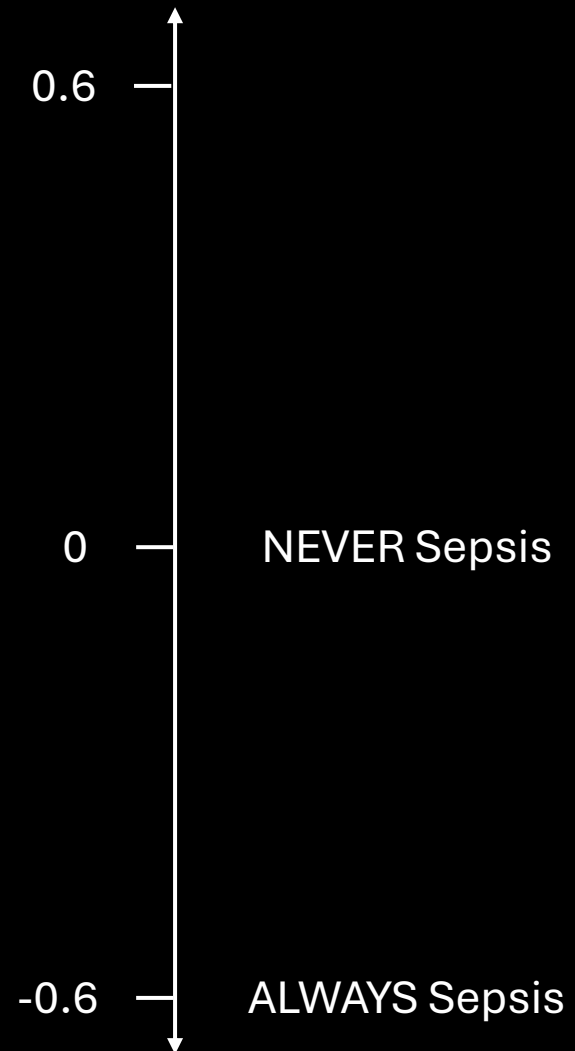
Goal: Optimise
for Clinical utility
Score

Utility Score



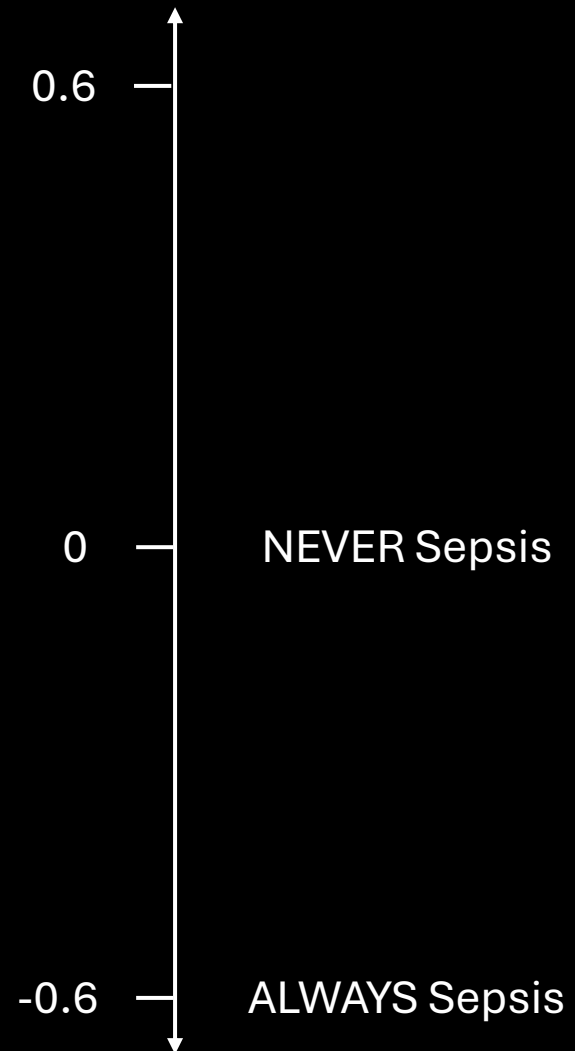
Model which ALWAYS predicts Sepsis
-> Reaches Utility of **-0.60**

Utility Score



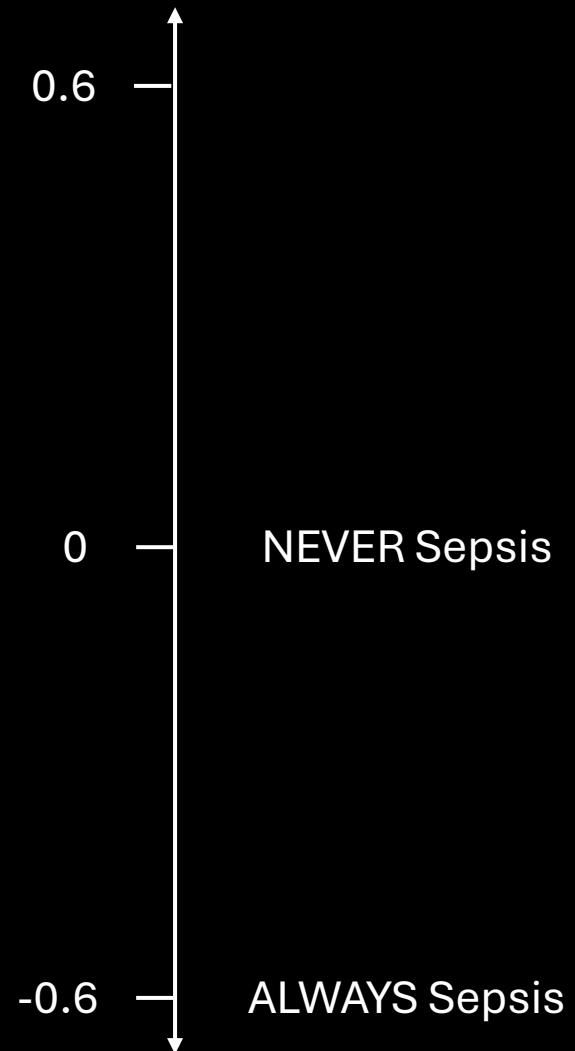
Model which NEVER predicts Sepsis
-> Reaches Utility of **0**

Utility Score

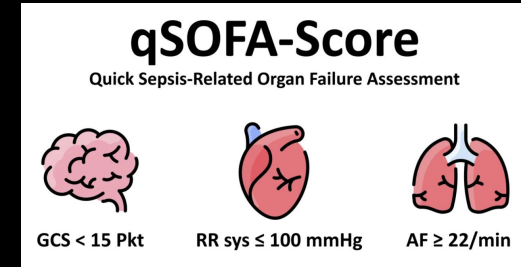


Model which NEVER predicts Sepsis
-> Reaches Utility of **0**

Utility Score



Utility Score

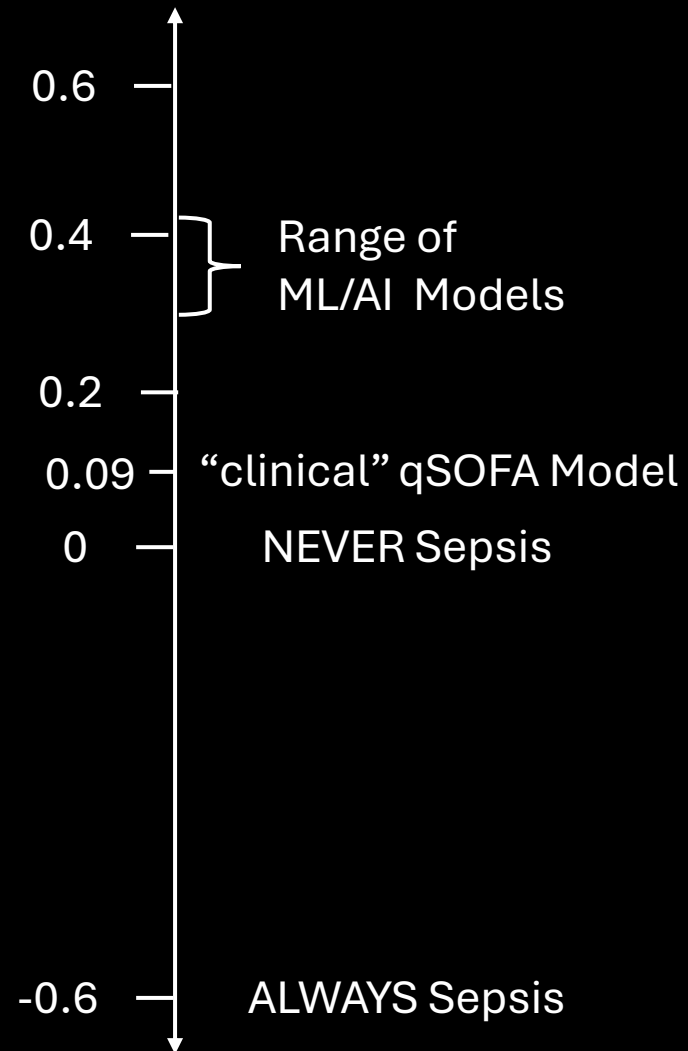


Model based on partial **qSOFA**
-> Reaches Utility of **0.09**

Utility Score



Utility Score



How do ML/AI
models perform?

Results



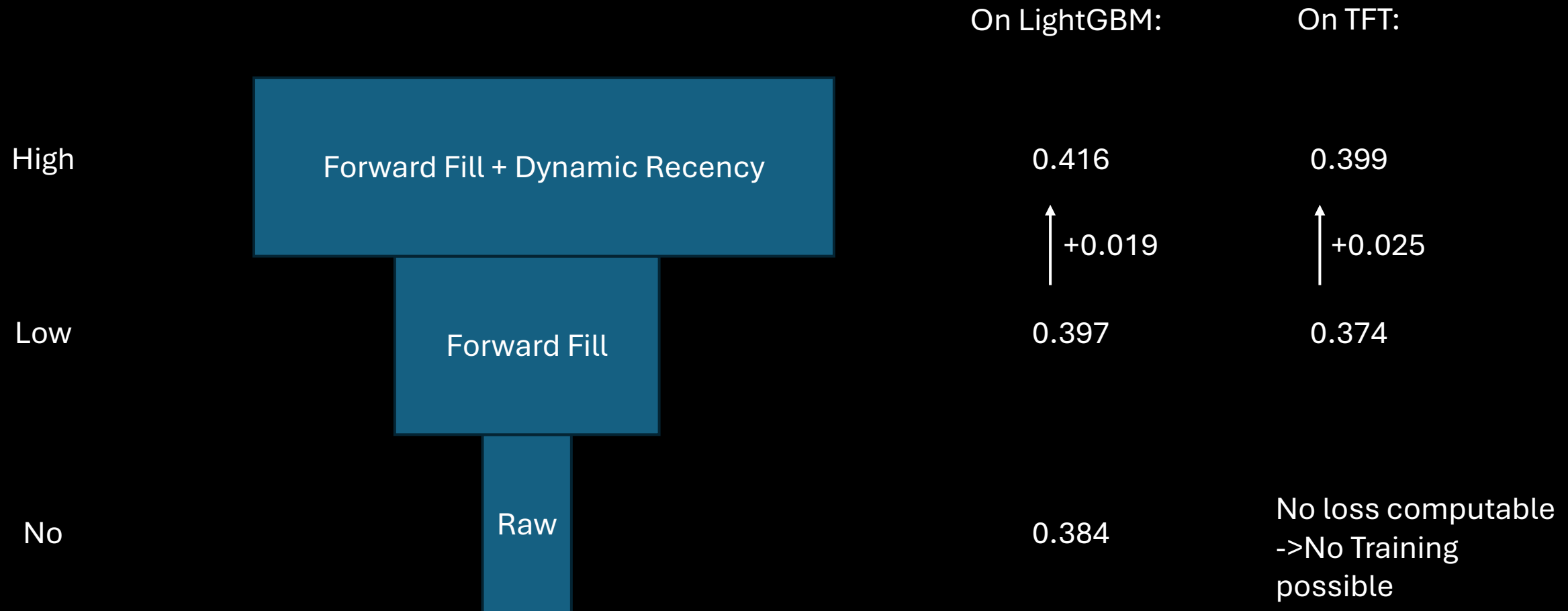
Other evaluated models included:

PatchTST 0.375

Crossformer 0.358

TimesNET 0.333

Results: Data Preprocessing Strategies



Limitations

- **Limited hyperparameter optimization & only single fold:**
Deep models were trained with minimal tuning (time/compute constraints) and on a single train/validation split (only LightGBM used 5-fold CV)
- **Regularization trade-offs:**
Overfitting was controlled using very strong dropout, gradient clipping and small learning rates. More data would likely improve generalization.

Discussion

LightGBM performs better when

- Data is limited
- Data is highly sparse (many missing values) and strongly imbalanced
- Fast training and robustness are prioritized

Temporal Fusion Transformer may perform better when

- Substantially more data is available
- Multimodal inputs are used (e.g., structured data combined with text or image embeddings)
- The task is more complex, (predicting multiple labels instead of a single label)
- Sufficient GPU compute is available

Future Work: Key Question

What are the **most important** improvements to make AI/ML models clinically useful?

Future Work

1

Explainable AI

Make TFT interpretable using built-in attention/feature-importance visualisations to better guide clinicians

2

Add clinical context

Incorporate clinical textual information (e.g. reason of hospital stay) and/or imaging (e.g. chest X-rays) via e.g. the TFT static metadata pathway.

3

Optimise for clinical utility

Train end-to-end on the utility score using Reinforcement Learning Approach (instead of probabilities + threshold tuning)

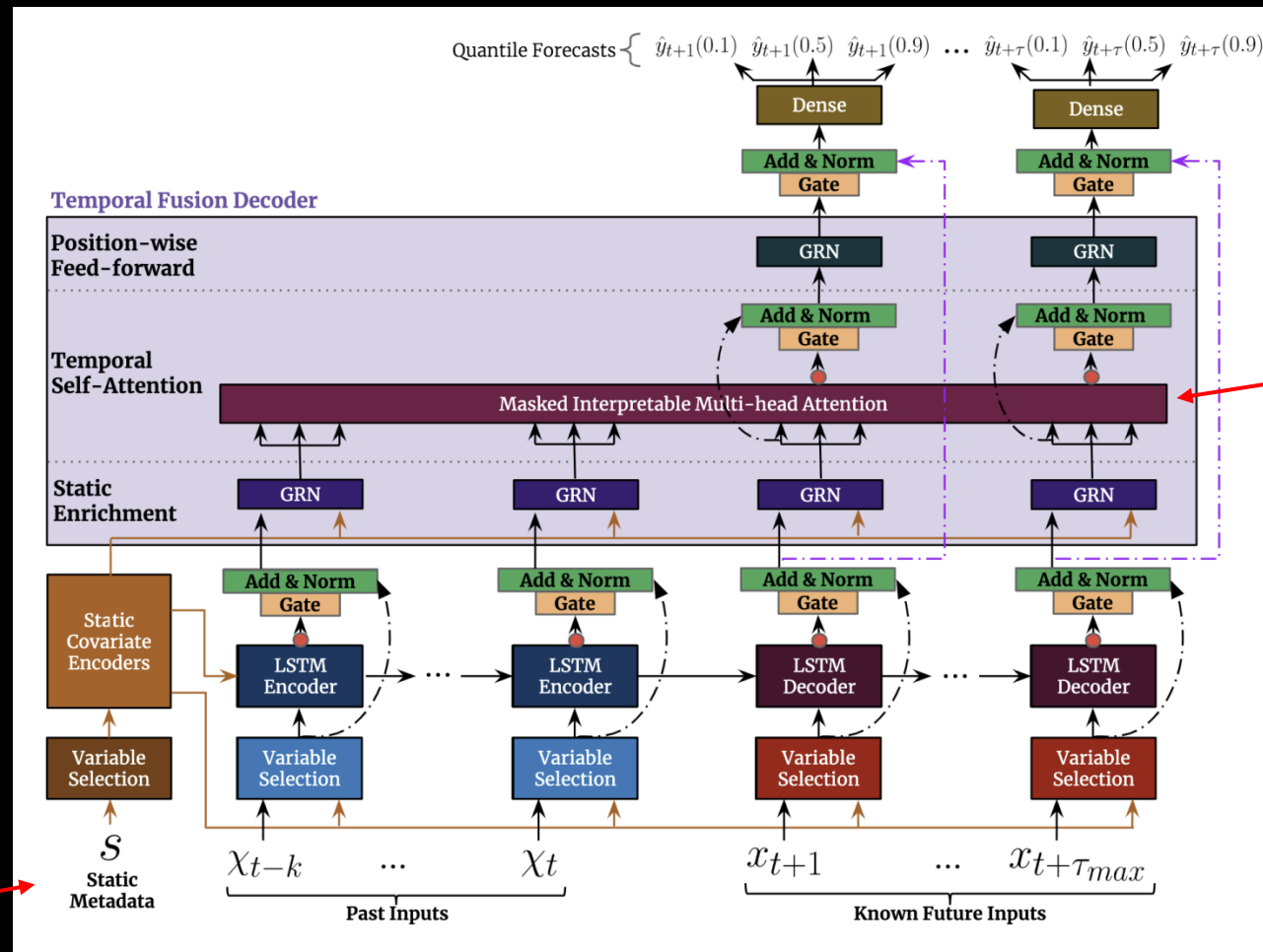
4

Scale data & tasks

Train on **larger datasets** like MIMIC-IV and/or eICU to reduce overfitting and help generalization and **add labels** (e.g., heart failure, ARDS) to improve accuracy.

Model Architecture: Temporal Fusion Transformers

Add clinical
context
(image or text
embeddings)



To conclude...

Using AI Prediction
Models have the
potential to save
lives!

Thank you for your attention!



GitHub Link to this project:
https://github.com/LucasNebelung/predicting_sepsis

