

Building Interpretable Models for Time Series Classification

Hugues Turbé, Mina Bjelogrić, Gianmarco Mengaldo, Medhi Namdar, Christian Lovis

MOTIVATIONS

Post-hoc Interpretability: Aim to identify features' contribution to model's outputs.

Interpretability methods attribute relevance to input features:

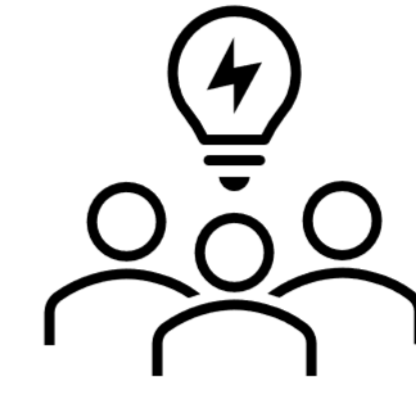
- Positive: input feature plays in favor of the model's prediction
- Negative: input feature plays against the model's prediction.

Interpretability is key for:

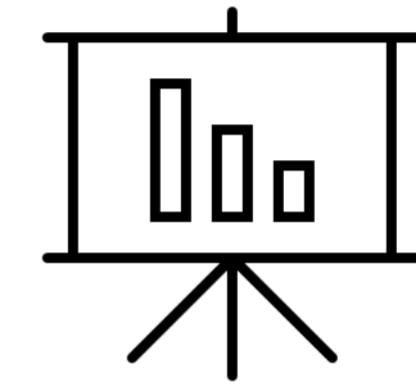
- **Building trust** in Artificial Intelligence (AI) based models¹
- **Regulatory aspect**².

CHALLENGES

- Provide **meaningful insights** into the model

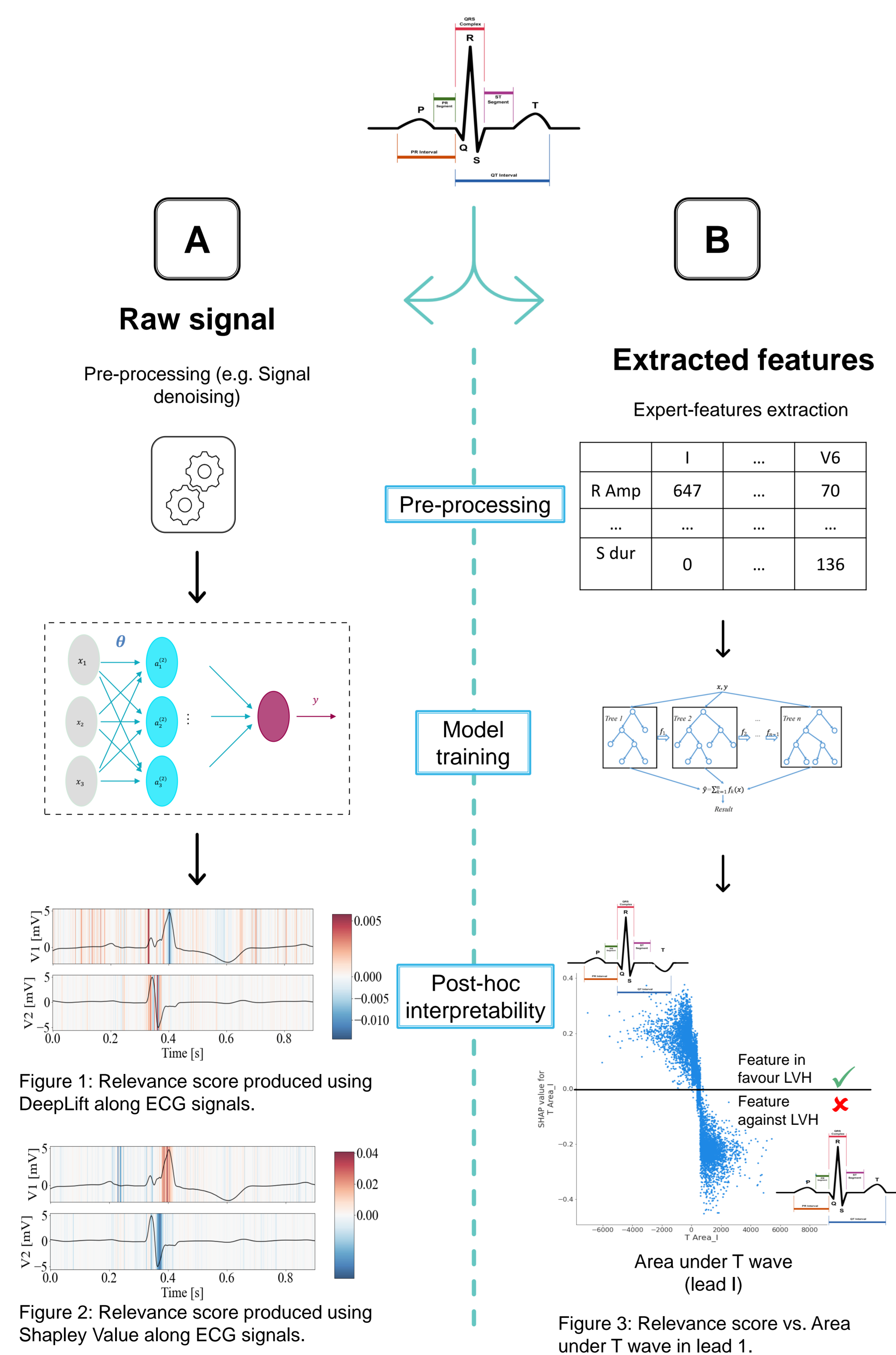


- **Evaluate** interpretability methods



CLASSIFICATION APPROACHES

Insights into the model is dictated by model's inputs.



METHODS

A

Classification Task³:

Right Bundle Branch Block (RBBB)

Models:

LSTM, CNN, Transformer (presented)

Interpretability methods:

DeepLift, DeepLiftShap, Shapley Value, KernelShap, GradShap

B

Classification Task⁴:

Left Ventricular Hypertrophy (LVH)

Model:

LightGBM

Interpretability methods:

Tree explainer

DISCUSSION

Different **interpretability methods** produce remarkably **different results³**:

- **Evaluation** of interpretability methods is **critical**.
- **Evaluation** should **assess the truthfulness** of the relevance attribution.

Relevance allows to compare features used by the model and cardiologists:

A: Shapley Value correctly identifies the R wave in lead V1 for LBBB classification.

B: Inversed T wave in lead I correctly identified as discriminator for LVH classification.

→ Build trust in models' predictions.



1. Shad, Rohan, et al. "Designing clinically translatable artificial intelligence systems for high-dimensional medical imaging." Nature Machine Intelligence 3.11 (2021): 929-935.
 2. European Commission, Directorate-General for Communications Networks, Content and Technology. Proposal for a regulation of the European parliament and of the council laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts (2021).
 3. Turbé, Hugues, et al. "InterpretTime: a new approach for the systematic evaluation of neural-network interpretability in time series classification." arXiv preprint arXiv:2202.05656 (2022).
 4. Turbé, Hugues, et al. "A Lightweight and Interpretable Model to Classify Bundle Branch Blocks from ECG Signals." Challenges of Trustable AI and Added-Value on Health. IOS Press, 2022. 43-47.