# 5. **Time-dependent explanations of neural networks for survival analysis

**Keywords:** tabular data, neural networks, experiments, pytorch

**Goal:** Compare SurvSHAP(t) model-agnostic explanation for survival models to explanations specific to neural networks, e.g. DeepLift. *Hopefully* model-specific explanations are comparable to SurvSHAP(t), but a lot faster to compute.
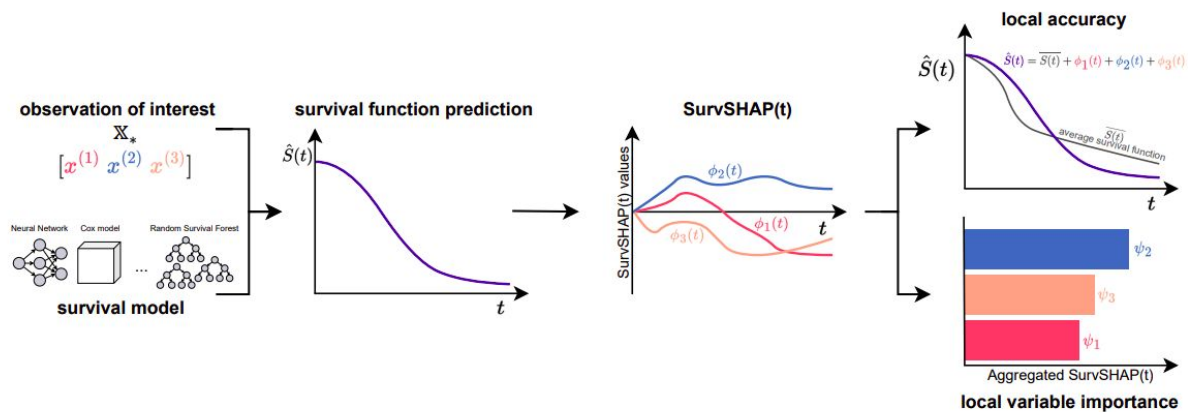
## Measure:

1. Computation time
2. Estimation error, i.e. how close is the faster method to SurvSHAP(t)
3. Evaluation measures, e.g. Avg-Sensitivity, Faithfulness Correlation

# SURVIVAL ANALYSIS

| | x0 | x1 | x2 | x3 | x4 | x5 | x6 | x7 | x8 | duration | event |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 5.603834 | 7.811392 | 10.797988 | 5.967607 | 1.0 | 1.0 | 0.0 | 1.0 | 56.840000 | 99.333336 | 0 |
| 1 | 5.284882 | 9.581043 | 10.204620 | 5.664970 | 1.0 | 0.0 | 0.0 | 1.0 | 85.940002 | 95.733330 | 1 |
| 3 | 6.654017 | 5.341846 | 8.646379 | 5.655888 | 0.0 | 0.0 | 0.0 | 0.0 | 66.910004 | 239.300003 | 0 |
| 4 | 5.456747 | 5.339741 | 10.555724 | 6.008429 | 1.0 | 0.0 | 0.0 | 1.0 | 67.849998 | 56.933334 | 1 |
| 5 | 5.425826 | 6.331182 | 10.455145 | 5.749053 | 1.0 | 1.0 | 0.0 | 1.0 | 70.519997 | 123.533333 | 0 |

# SurvSHAP(t)

# DeepHit



Figure 2: The architecture of DeepHit with two competing events.

# DeepLIFT

# Integrated Gradients

# Metrics

**On the (In)fidelity and Sensitivity of Explanations**

**Evaluating and Aggregating Feature-based Model Explanations**