# Learning the Optimal Stopping for Early Classification within Finite Horizons via Sequential Probability Ratio Test

**Akinori F. Ebihara (aebihara@nec.com),**   Taiki Miyagawa,   Kazuyuki Sakurai,   Hitoshi Imaoka      **NEC Corporation, Japan**
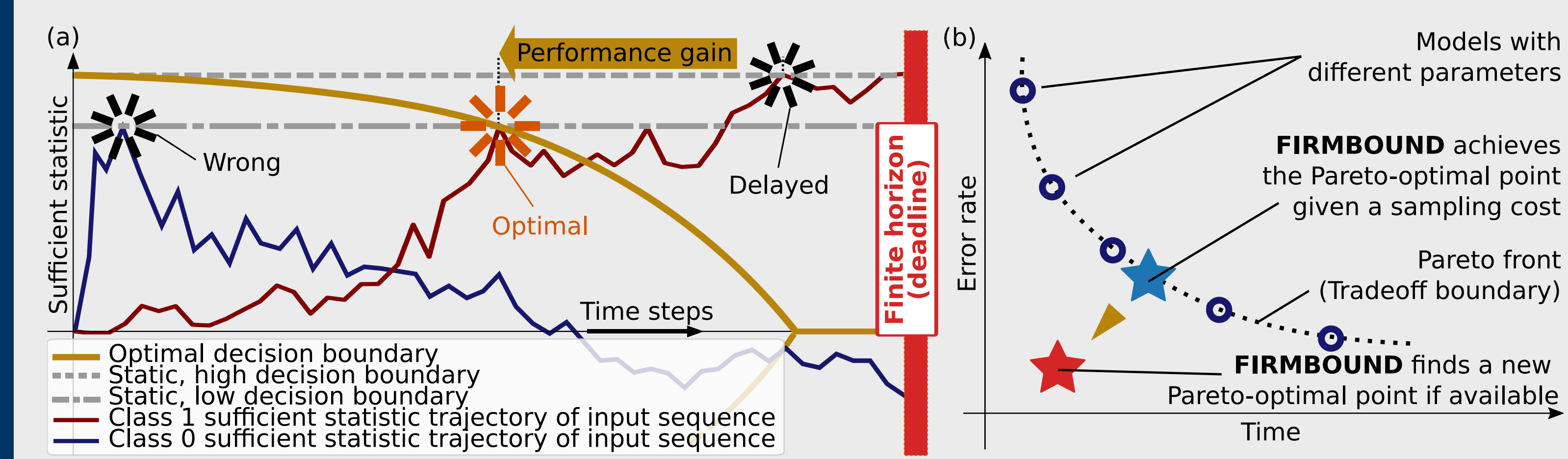
## Early Classification Is Pivotal for Time-Sensitive Analysis



Evidence   Left   Right   Biometric entrance gate

## Sequential Probability Ratio Test (SPRT) Is Suboptimal Under *Finite Horizon* Constraints



(a) Performance gain / Wrong / Optimal / Delayed / Time steps / Finite horizon (deadline)

Optimal decision boundary
Static, high decision boundary
Static, low decision boundary
Class 1 sufficient statistic trajectory of input sequence
Class 0 sufficient statistic trajectory of input sequence $\mathscr{S}_t$

(b) Models with different parameters / **FIRMBOUND** achieves the Pareto-optimal point given a sampling cost / Pareto front (Tradeoff boundary) / **FIRMBOUND** finds a new Pareto-optimal point if available

## Finding Optimal Thresholds Is Computationally Intensive

**<u>Bayes Risk</u>** to be minimized with the optimal threshold

Sufficient statistic i.e. Log-likelihood ratio (LLR) or posterior   Decision to classify as class k   Sequential data   Class posterior   Constant   time

$$\text{APR}_t(\mathscr{S}_t, d_t(X^{(1,t)}) = k) := \bar{L}_k(1 - \pi_k(X^{(1,t)})) + ct$$

A posteriori risk (APR)   Constant   Misclassification Penalty   Sampling cost

**<u>Backward induction equation</u>** to minimize the Bayes Risk

Curse of dimensionality at estimating the conditional expectation

Continuation risk
$$\tilde{G}_t(\mathscr{S}_t) = \mathbb{E}\left[G_{t+1}^{\min}(\mathscr{S}_{t+1})|\mathscr{S}_t\right] + c$$

Stopping risk
$$G_t^{st}(\mathscr{S}_t) = \min_k\left\{\bar{L}_k(1 - \pi_k(X^{(1,t)}))\right\}, \quad \text{Time } t \in \{1, \cdots, T\}$$

Minimum risk
$$G_t^{\min}(\mathscr{S}_t) := \begin{cases} G^{st}(\mathscr{S}_t) & (t = T) \\ \min\left\{G^{st}(\mathscr{S}_t), \tilde{G}_t(\mathscr{S}_t)\right\} & (1 \le t < T). \end{cases}$$

## FIRMBOUND Is "Doubly-Consistent"
- providing consistent estimation of likelihood ratio and the thresholds

**Noisy convex regression**

Total #data   Concave function   Continuation risk with noise   Regularization term   Sampling cost

$$\hat{\hat{G}}_t(\{X_m^{(1,T)}\}_{m=1}^M) \in \arg\min_{f: \text{concave}}\left\{\frac{1}{M}\sum_{m=1}^{\hat{M}}\left(f(\mathscr{S}_t(X_m^{(1,t)})) - \mathscr{G}_m^{(t)}\right)^2 + \lambda\|f\|\right\} + c$$

Mean squared error   Constant

Continuation risk
$$\mathscr{G}_m^{(t+1)} = \tilde{G}_t(\mathscr{S}_t(X_m^{(1,t)})) + \epsilon_m^{(t)} - c$$

Sequential data   Noise

cf) Siahkamari et al. 2022

**(Optional) Gaussian Process (GP) regression for lightweight training**

- Significantly reduces training time with potential compromise in statistical consistency
- Assuming $\epsilon_m^{(t)}$ and $\{\tilde{G}_t(\mathscr{S}_{t,m})\}_{m\in[M]}$ are Gaussian noise and GP, respectively, the conditional expectation estimation problem is reduced down to GP regression problem:

$$\mathscr{G}_m^{(t+1)} + c = \tilde{G}_t(\mathscr{S}_{t,m}) + \epsilon_m^{(t)}$$

Response variable   Latent function   Gaussian noise   Explanatory variable

cf) Ebihara et al. 2021, Ebihara and Miyagawa, 2021

**Likelihood ratio estimation with SPRT-TANDEM**

Trainable parameters   Ground-truth label   Total #classes   Total #data   Total #data per class   Estimated LLR

$$\hat{L}_{\text{LSEL}}(\boldsymbol{w}; \{(X_m^{(1,T)}, y_m\}_{m\in[M]}) := \frac{1}{KM}\sum_{k\in[K]}\sum_{t\in T}\frac{1}{M_k}\sum_{i\in I_k}\log(1 + \sum_{l(\ne k)\in[K]}e^{-\hat{\lambda}_{kl}(\boldsymbol{w}, X^{(1,t)})})$$

Loss function LSEL

$$\hat{\lambda}_{kl}(X^{(1,t)}) = \sum_{s=N+1}^t \log\frac{\pi_k(X^{(s-N,s)})}{\pi_l(X^{(s-N,s)})} - \sum_{s=N+2}^t \log\frac{\pi_k(X^{(s-N,s-1)})}{\pi_k(X^{(s-N,s-1)})} - \log\chi_{kl}$$
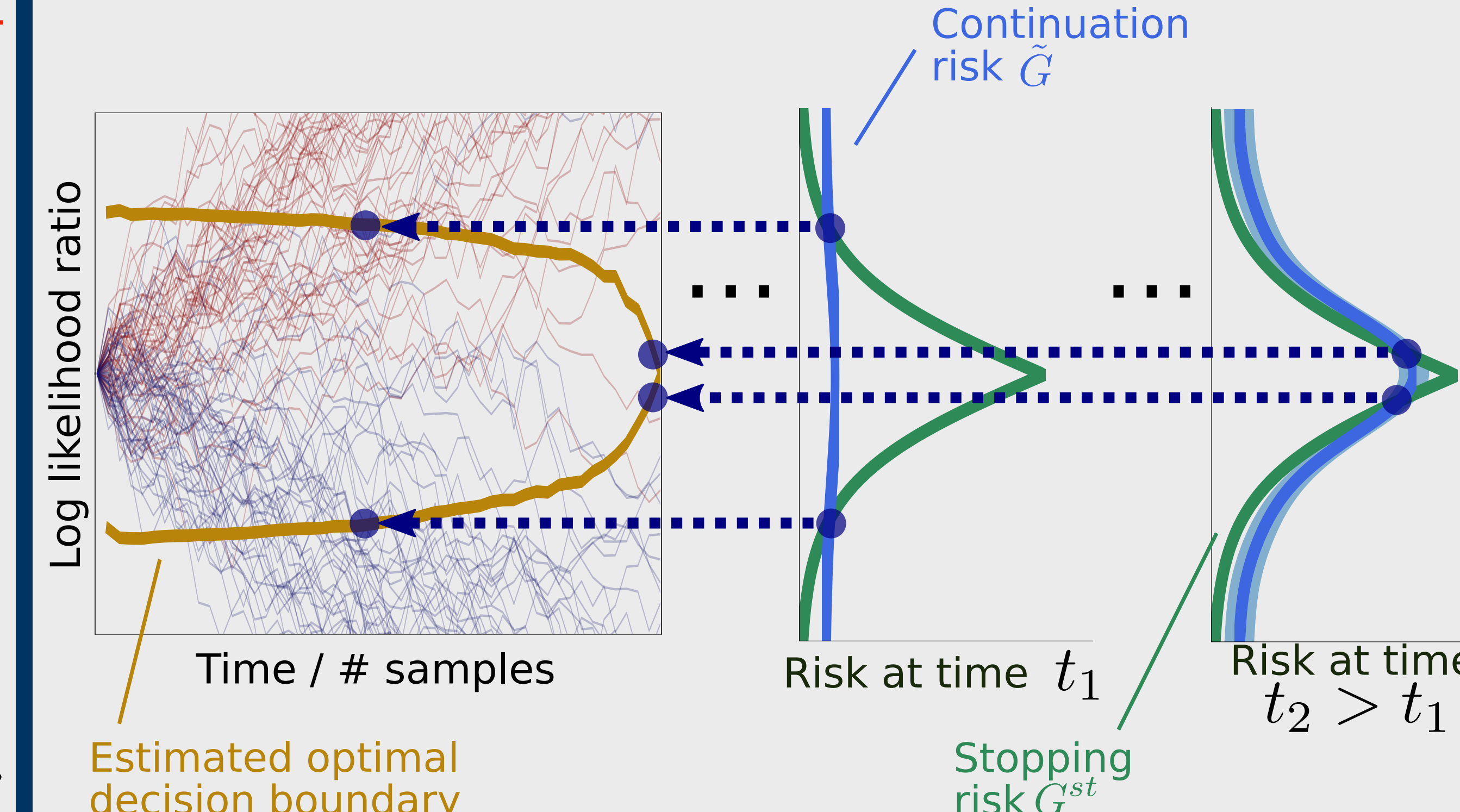
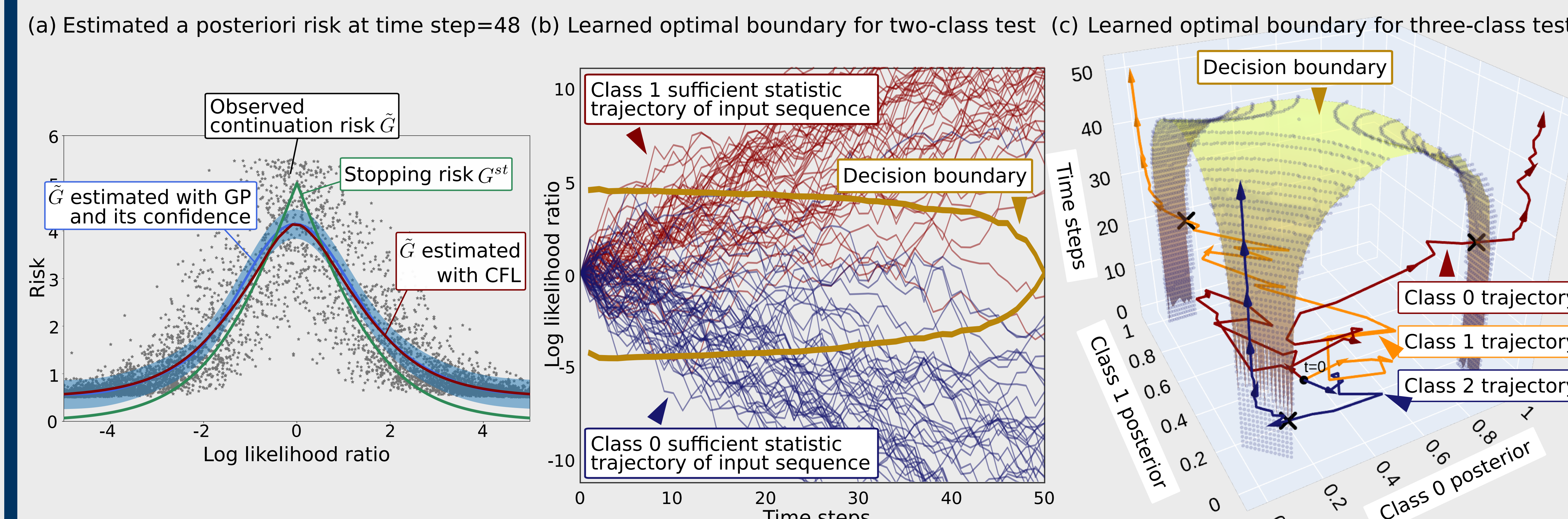LLR estimated with TANDEM formula   Class posterior   Order of Markov process   Class prior ratio

## Optimal Threshold Defined by the Intersection of Two Curves



Continuation risk $\tilde{G}$   Stopping risk $G^{st}$

Time / # samples   Risk at time $t_1$   Risk at time $t_2 > t_1$

Estimated optimal decision boundary

## Example Estimated Thresholds on Gaussian Datasets



(a) Estimated a posteriori risk at time step=48   (b) Learned optimal boundary for two-class test   (c) Learned optimal boundary for three-class test

Observed continuation risk $\tilde{G}$
$\tilde{G}$ estimated with GP and its confidence
$\tilde{G}$ estimated with CFL
Stopping risk $G^{st}$

Class 1 sufficient statistic trajectory of input sequence
Decision boundary
Class 0 sufficient statistic trajectory of input sequence

Class 0 trajectory / Class 1 trajectory / Class 2 trajectory

## Training Pipeline



Forward estimation

$$\lambda_t = \text{TANDEM}\left(\{X_m^{(i-N,i)}, X_m^{(j-N,j-1)}|i,j\in\mathbb{N}, N+1\le i\le t, N+2\le j\le t\}\right) \text{ if } N < t; \text{ else } X_m^{(1,t)}$$

Forward estimation of log likelihood ratio distributions

Conditional expectations estimated with **GP/CFL**

$$\hat{\mathbb{E}}\left[\min\left\{\tilde{G}_{T-1}, G_{T-1}^{st}\right\}|\lambda_{T-2}\right]$$

$$\hat{\mathbb{E}}\left[\min\left\{\tilde{G}_{T-2}, G_{T-2}^{st}\right\}|\lambda_{T-3}\right]$$

Backward induction

## FIRMBOUND Minimizes the Bayes Risk to Delineate the Pareto Front



(a) Two-class i.i.d. Gaussian   (b) Three-class i.i.d. Gaussian   (c) Damped oscillating non-i.i.d. LLRs (DOL)
(d) SiW   (e) HMDB51   (f) UCF101

Mean hitting time (#frames)

$c = 2L/T$ FIRMBOUND (CFL) on true $\mathscr{S}_t$
$c = L/T$ FIRMBOUND (CFL) on true $\mathscr{S}_t$
$c = 0.1L/T$ FIRMBOUND (CFL) on estimated $\mathscr{S}_t$
$c = 2L/T$ FIRMBOUND (CFL) on estimated $\mathscr{S}_t$
$c = L/T$ FIRMBOUND (CFL) on estimated $\mathscr{S}_t$
$c = 0.1L/T$ FIRMBOUND (GP) on estimated $\mathscr{S}_t$
$c = 2L/T$ Vanilla SPRT on true LLRs
$c = L/T$ Vanilla SPRT on true LLRs
$c = 0.1L/T$ Vanilla SPRT on estimated LLRs (SPRT-TANDEM)

Static SPRT on true LLRs
Static SPRT on estimated LLRs (SPRT-TANDEM)

LSTMms / EARLIEST / TCNTransformer / CALIMERA

| Dataset | Gauss2est. | Gauss3est. | DOL | SiW | HMDB | UCF101 | FordA |
|---|---|---|---|---|---|---|---|
| Trial repeats | 5 | 3 | 3 | 5 | 6 | 10 | 2 |
| ↓MVHT, vanilla SPRT with static threshold | 10.47 | 44.02 | 489.89 | 2.87 | 199.31 | 0.55 | 32.15 |
| ↓MVHT, FIRMBOUND with CFL | 9.01 | 42.78 | 405.78 | 1.97 | 195.35 | 0.53 | 23.39 |
| ↑Difference in MVHT (positive is better) | 1.45 | 1.24 | 84.11 | 0.90 | 3.97 | 0.017 | 8.76 |

## TL;DR

**<u>FIRMBOUND</u>** is an **<u>SPRT-based early classification</u>** framework that provides a **<u>statistically consistent and computationally efficient estimator of optimal decision boundaries</u>** for time series of finite lengths, tailored for large-scale real-world problems.