



# BÁO CÁO DỰ ÁN

## ĐỀ TÀI:

**SAGE (Sentiment Analysis & Gap Execution)**  
– Biến tin tức tiếng Việt thành tín hiệu rủi  
ro thị trường và giao dịch theo thời gian thực

**Đội thi: AI Max – Trường Đại học Công nghệ**  
– **ĐHQGHN**

Hà Nội, năm 2025

# Tóm tắt

**SAGE** là hệ thống AI thời gian thực chuyển hóa tin tức tiếng Việt thành thước đo rủi ro và tín hiệu giao dịch cho thị trường cổ phiếu. Hệ thống gồm các bước: suy luận cảm xúc bằng **PhoBERT**, trích xuất đặc trưng tin tức, dự báo chế độ biến động/jump (mô hình A), dự báo small-gap (mô hình B), và kiểm soát rủi ro bằng logic cổng.

Hệ thống được triển khai bằng **FastAPI** với hàng đợi inference, ghi log và giám sát drift. Kết quả cho thấy mô hình A phân biệt jump khá tốt (AUROC  $\approx 0.70$ ), mô hình B tạo tín hiệu small-gap thận trọng và mang lại hiệu quả rủi ro điều chỉnh dương. **SentixFlow** hiện ở mức nguyên mẫu hoàn thiện, hướng đến mở rộng thích ứng theo phiên và đánh giá trực tiếp trên OMS.

## 1 Giới thiệu

Thị trường Việt Nam có tỷ trọng nhà đầu tư cá nhân cao nên tin tức tác động mạnh đến kỳ vọng và biến động ngắn hạn. Tuy vậy, việc biến tín hiệu tin tức thành quyết định giao dịch khả thi gặp các thách thức: (i) tín hiệu suy giảm nhanh theo thời gian; (ii) rủi ro rò rỉ thông tin khi gán nhãn; (iii) dịch chuyển phân phối (drift) của đặc trưng theo giai đoạn; và (iv) phụ thuộc chế độ thị trường (trending vs. mean-reverting), nhất là quanh *gap* sau mở cửa.

Chúng tôi đề xuất **SentixFlow**—một *pipeline* AI thời gian thực cho thị trường Việt Nam, gồm hai lớp mô hình bổ sung:

- **Model A (Vol/Jump)**: dự báo chế độ biến động (*regime*) và xác suất *jump*, dùng để quy mô vị thế theo ánh xạ  $\{low : 1.0, mid : 0.7, high : 0.4\}$ .
- **Model B (Small-gap)**: phân loại *fade/follow* và hồi quy **ret\_OC** cho các phiên *gap* nhỏ, kèm *gates* (ngưỡng *sentiment/entropy*, *fade-rate* cuốn trượt, *edge* tối thiểu sau phí) để giảm nhiễu.

Đặc trưng tin tức được rút trích bằng **PhoBERT** (wonrax/phobert-base-vietnamese-sentiment) và chuẩn hoá z-score/cửa sổ lăn. Nhãn mục tiêu (*gap*, **ret\_OC**, *realized volatility*, *regime*) được xây bằng phân vị mở rộng theo thời gian nhằm loại bỏ *look-ahead*. Hệ thống *serving* (**FastAPI**, GPU/CPU) ghi log JSONL theo ngày; mô-đun *monitoring* theo dõi *drift* (PSI), phân phối đặc trưng, độ trễ/thông lượng và tỷ lệ ngày bị chặn bởi *gates*. Kết quả xuất ra *risk gauge* (*regime*, *jump*), *trading signal* (*fade/follow/skip*, *size*, *audit*) và báo cáo giám sát phục vụ tích hợp OMS/Telegram/Slack.

Hệ thống đạt mức nguyên mẫu hoàn thiện: *end-to-end* từ tin tức  $\rightarrow$  đặc trưng  $\rightarrow$  rủi ro  $\rightarrow$  tín hiệu, kèm *audit* và *monitoring*. Hạn chế hiện tại là dữ

liệu cấp chỉ số/ngày và chưa phân giải theo mã/ngành; hướng phát triển gồm *calibration* theo phiên, *domain adaptation* theo nguồn tin, thêm đặc trưng liên thị trường và mở rộng sang nhiều lớp tài sản.

## 2 Phương pháp

### 2.1 Bài toán và ký hiệu

Ký hiệu  $t$  là ngày giao dịch trên thị trường Việt Nam (UTC+7). Với mỗi ngày  $t$ , ta quan sát tập tin tức  $\mathcal{N}_t = \{n_{t,i}\}_{i=1}^{M_t}$  và dữ liệu thị trường của chỉ số đại diện VN30 gồm giá mở cửa  $O_t$ , cao nhất  $H_t$ , thấp nhất  $L_t$ , đóng cửa  $C_t$ . Suất sinh lợi được định nghĩa:

$$r_{t+1}^{\text{gap}} = \frac{O_{t+1} - C_t}{C_t}, \quad r_{t+1}^{\text{oc}} = \frac{C_{t+1} - O_{t+1}}{O_{t+1}}, \quad r_{t+1}^{\text{1d}} = \frac{C_{t+1} - C_t}{C_t}.$$

Ràng buộc kiểm tra:  $(1 + r_{t+1}^{\text{1d}}) \approx (1 + r_{t+1}^{\text{gap}})(1 + r_{t+1}^{\text{oc}})$ .

**Độ biến động thực tế và chế độ.** Độ biến động ngày kế tiếp được ước lượng theo công thức Parkinson (high–low):

$$\text{RV}_{t+1}^{\text{PK}} = \sqrt{\frac{1}{4 \ln 2} \left( \ln \frac{H_{t+1}}{L_{t+1}} \right)^2}, \quad \text{RV}_{t+1}^{(5)} = \text{stdev}(r_{t-k}^{\text{1d}})_{k=0}^4.$$

Để tránh *look-ahead*, chế độ biến động được xác định bằng các phân vị *expanding* chỉ dùng dữ liệu tới thời điểm  $t$ :

$$q_{p,t} = \text{Quantile}_p(\{\text{RV}_\tau^{(5)}\}_{\tau \leq t}), \quad \text{min\_periods} = 60.$$

Nhãn cho ngày  $t+1$ :

$$R_{t+1} = \{0 \text{ nuRV}_{t+1}^{(5)} \leq q_{33,t}, 1 \text{ nu} q_{33,t} < \text{RV}_{t+1}^{(5)} < q_{67,t}, 2 \text{ nu} \text{RV}_{t+1}^{(5)} \geq q_{67,t}\}.$$

**Phân nhóm gap và nhãn fade.** Ta phân nhóm độ lớn tuyệt đối của gap tại  $t+1$ :

$$B_{t+1} = \{ \text{micronu} |r_{t+1}^{\text{gap}}| < 0.3\%, \text{smallnu} 0.3\% \leq |r_{t+1}^{\text{gap}}| < 1.0\%, \text{largenu} |r_{t+1}^{\text{gap}}| \geq 1.0\% \}.$$

Mục tiêu *fade* được định nghĩa:

$$\text{Fade}_{t+1} = 1\{\text{sign}(r_{t+1}^{\text{gap}}) \neq \text{sign}(r_{t+1}^{\text{oc}})\}.$$

**Nhân jump.** Một ngày được coi là *jump* nếu suất sinh lợi tuyệt đối vượt quá ngưỡng nhiều lần độ lệch chuẩn trượt:

$$\text{Jump}_{t+1} = 1 \left\{ |r_{t+1}^{\text{1d}}| > \tau \cdot \sigma_t \right\}, \quad \sigma_t = \text{stdev}(r_{t-k}^{\text{1d}})_{k=1}^{20}, \quad \tau = 2.$$

## 2.2 Suy luận tin tức và tổng hợp hằng ngày

**Sentiment PhoBERT.** Mỗi bài báo  $n_{t,i}$  được chấm điểm bởi bộ phân loại cảm xúc tiếng Việt (PhoBERT, checkpoint `wonrax/phobert-base-vietnamese-sentiment`) với xác suất  $p_{t,i}^{\text{NEG}}, p_{t,i}^{\text{NEU}}, p_{t,i}^{\text{POS}}$ . Ta rút trích:

$$s_{t,i} = p_{t,i}^{\text{POS}} - p_{t,i}^{\text{NEG}}, \quad H_{t,i} = - \sum_{k \in \{\text{NEG}, \text{NEU}, \text{POS}\}} p_{t,i}^k \log p_{t,i}^k.$$

Đặc trưng theo ngày:

$$\bar{s}_t = \frac{1}{M_t} \sum_i s_{t,i}, \quad \tilde{s}_t = \text{median}_i(s_{t,i}), \quad \text{std}_s(t) = \text{stdev}_i(s_{t,i}),$$

$$\bar{H}_t = \frac{1}{M_t} \sum_i H_{t,i}, \quad \text{std}_H(t) = \text{stdev}_i(H_{t,i}), \quad M_t = \text{news\_count}(t).$$

**Chuẩn hoá  $z$ -score và spike.** Với cửa sổ mở rộng hoặc trượt  $\mathcal{W}_t$  (ví dụ 120 ngày) ta tính:

$$z_x(t) = \frac{x_t - \mu_x(t)}{\sigma_x(t)}, \quad z_x^\lambda(t) = \frac{x_t - \mu_x^\lambda(t)}{\sigma_x^\lambda(t)}$$

trong đó  $\mu_x^\lambda, \sigma_x^\lambda$  là các ước lượng suy giảm theo mũ (tham số quên  $\lambda \in (0, 1)$ , mặc định  $\lambda = 0.98$ ). Áp dụng winsorization tại  $\pm 5\sigma$ . Cờ spike:  $1\{z_M(t) > 1\}, 1\{z_M(t) > 2\}$ .

## 2.3 Mô hình: Hai tầng (A+B)

**Model A (Risk gauge).**

- **Hồi quy biến động:**  $f^{\text{rv}} : X_A(t) \mapsto \widehat{\text{RV}}_{t+1}^{\text{PK}}$  huấn luyện bằng HistGradientBoostingRegressor (HGB).
- **Phân loại chế độ:**  $g^{\text{reg}} : X_A(t) \mapsto \hat{R}_{t+1} \in \{0, 1, 2\}$  qua Logistic Regression đa thức (one-vs-rest).

- **Phân loại jump:**  $h^{\text{jump}} : X_A(t) \mapsto \hat{p}_{\text{jump}}(t+1)$  qua HGB với trọng số lớp ( $w_+$ ) và suy giảm gần đây  $\gamma$ .

$X_A(t)$  chứa đặc trưng tin tức tổng hợp:  $\{\bar{s}_t, \tilde{s}_t, \text{std}_s(t), \bar{H}_t, \text{std}_H(t), z\text{-score}, \text{spike}, \text{và các tóm tắt trượt ngắn hạn (ví dụ trung bình/cực trị 5 ngày)}\}$ . Quy mô vị thế được ánh xạ:

$$S(\hat{R}) = \{1.0, 0.7, 0.4\} \quad \text{ngvi} \hat{R} \in \{0, 1, 2\}.$$

**Model B (Quyết định small-gap).** Với các *gap nhỏ* ( $B_{t+1} = \text{small}$ ), ta huấn luyện: Bộ phân loại  $c : X_B(t) \mapsto \hat{p}_{\text{fade}}(t+1)$ ,  $Bhiquyr : X_B(t) \mapsto \hat{r}_{t+1}^{\text{oc}}$ . Dùng Logistic Regression cho  $c$  và Random Forest (hoặc HGB) cho  $r$ . Ngưỡng xác suất  $\theta$  được chọn trên tập validation để tối ưu F1 hoặc một hàm mục tiêu thân thiện với AP.  $X_B(t)$  kế thừa đặc trưng tin tức và thêm dấu/hệ số của gap hiện tại.

## 2.4 Logic quyết định và các gates

Với đặc trưng ngày  $t$  và gap quan sát  $g = r_{t+1}^{\text{gap}}$  sau mở cửa:

1. Tính rủi ro:  $\hat{R} = g^{\text{reg}}(X_A)$ ,  $\hat{p}_{\text{jump}} = h^{\text{jump}}(X_A)$ , quy mô  $S$ .
2. Nếu  $\hat{p}_{\text{jump}} > \eta$  thì giảm/skip theo chính sách ( $\eta$  chọn trên validation).
3. Nếu  $|g|$  thuộc nhóm *small*, đánh giá  $c$  và  $r$ :

$$\hat{p}_{\text{fade}} = c(X_B), \quad \hat{r}^{\text{oc}} = r(X_B).$$

4. Áp dụng các *gates* (kiểm soát chất lượng):
  - (G1)  $|\bar{s}_t| \leq s_{\text{max}}$ ;
  - (G2)  $z_M^\lambda(t) \leq z_{\text{news}}^{\text{max}}$ ;
  - (G3)  $z_{\text{std}_H}^\lambda(t) \leq z_{\text{ent}}^{\text{max}}$ ;
  - (G4)  $\hat{p}_{\text{fade}} \geq \theta$ ;
  - (G5)  $\hat{r}^{\text{oc}} \leq r_{\text{fade}}^{\text{max}}$  (với fade),  $\hat{r}^{\text{oc}} \geq r_{\text{follow}}^{\text{min}}$  (với follow);
  - (G6)  $\hat{F}_{30}(t) \geq F^{\text{min}}$  (tỷ lệ fade trượt, chế độ contrarian).
5. Quyết định:

$$\text{Nugatesfail} \Rightarrow \text{skip}; \quad \text{nusign}(g) \cdot \hat{r}^{\text{oc}} < 0 \vee \hat{p}_{\text{fade}} \geq \theta \Rightarrow \text{fade};$$

$$\text{nur}^{\hat{r}^{\text{oc}}} \geq r_{\text{follow}}^{\text{min}} \Rightarrow \text{follow}; \quad \text{ngclis} \Rightarrow \text{skip}.$$

Kích thước giao dịch =  $S(\hat{R})$ .

## 2.5 Huấn luyện, validation và kiểm thử

Chia tách theo thời gian, bám sát cấu trúc thị trường:

- **Huấn luyện:** đến 31/12/2024,
- **Validation:** 02/01/2025 đến 31/03/2025,
- **Kiểm thử:** 01/04/2025 đến 30/06/2025.

Thước đo: RMSE/MAE cho  $RV^{PK}$ , accuracy/F1 cho  $R$ , AUROC/AP cho  $Jump$ , AUROC/AP/F1 cho fade classifier, correlation/RMSE/MAE cho  $r^{oc}$ .

## 2.6 Triển khai và tổng hợp thời gian thực

**Dịch vụ sentiment.** PhoBERT chạy sau FastAPI (GPU/CPU) với hàng đợi yêu cầu. Endpoint `/infer` nhận JSON ...

# 3 Thực nghiệm

## 3.1 Dữ liệu & thu thập

**Nguồn tin tức.** Chúng tôi *tự thu thập* dữ liệu bài viết tiếng Việt từ CafeF, lưu thành bảng ngày `news_clean.parquet` (2015-01-08 → 2025-06-30, 3,212 ngày). Mỗi bản ghi gồm ngày, tiêu đề/nội dung đã làm sạch, và các trường kỹ thuật (ID, nguồn).

**Thị trường.** Dựa trên rổ VN30, chúng tôi xây dựng ba proxy (EW/VW/LW) và các nhãn mục tiêu ở Phase 1-3 (gap  $r_{t+1}^{gap}$ , return trong ngày  $r_{t+1}^{oc}$ , biến động thực  $RV_{t+1}^{PK}$ , regime 3 mức, jump). Tập ghép đầy đủ dùng cho học máy (`dataset_ab.clean.parquet`) có 1,088 ngày (02/2021 → 06/2025; scheme `vw` trong thí nghiệm).

## 3.2 Tiền xử lý & gắn nhãn

**Kiểm tra đồng nhất.** Bảo toàn đẳng thức  $(1 + r^{1d}) \approx (1 + r^{gap})(1 + r^{oc})$  với sai số trung bình  $\leq 2 \cdot 10^{-5}$ .

**Regime không rò rỉ.** Thay vì quantile toàn cục, dùng *expanding quantiles* (min\_periods=60) để gắn `regime3_t1`. Phân phối cuối: Low  $\approx 52\%$ , Mid  $\approx 30\%$ , High  $\approx 18-19\%$ .

**Jump.** Nhãn `jump_t1` khi  $|r^{1d}| > \tau \cdot \text{stdev}_{20}$  (mặc định  $\tau=2$ ); tỷ lệ jump  $\approx 6.5\%-7\%$ .

**Buckets gap.**  $|\text{gap}| < 0.3\%$  (micro),  $[0.3, 1)\%$  (small),  $\geq 1\%$  (large). Tương

quan  $\text{corr}(\text{gap}, r^{\text{oc}}) \approx -0.08 \sim -0.11$  (small-gap thiên về *fade*).

**Đặc trưng tin tức.** Mỗi bài được suy luận bằng PhoBERT (`wonrax/phobert-base-vietnamese-` cho  $p_{\text{NEG}}, p_{\text{NEU}}, p_{\text{POS}}$ ; tính  $s = p_{\text{POS}} - p_{\text{NEG}}$  và entropy  $H$ . Gộp ngày: trung bình/median/độ lệch chuẩn của  $s$  &  $H$ , chuẩn hoá  $z$  mở rộng, phiên bản làm tròn mũ ( $z_w$ ), cờ spike ( $z > 1, 2$ ). Các cột khoá ở Phase 6: `sent_score_mean`, `sent_entropy_std_z_w`, `news_count_z_w`.

### 3.3 Thiết lập huấn luyện

Chia thời gian: **Train**  $\leq$  2024-12-31, **Val** 2025-01-02  $\rightarrow$  2025-03-31, **Test** 2025-04-01  $\rightarrow$  2025-06-30.

**Model A (Risk gauge).** HGB cho hồi quy  $RV^{\text{PK}}$ ; Logistic Regression đa lớp cho **regime3**; HGB cho **jump** (class weight dương, trọng số phân rã gần trên mẫu).

**Model B (quyết định small-gap).** Logistic Regression cho  $p(\text{fade})$  và Random Forest (hoặc HGB) cho dự báo  $r^{\text{oc}}$ . Ngưỡng xác suất tối ưu theo F1/AP trên **val**.

**Chi phí giao dịch.** Mặc định 5bps. Mọi đặc trưng/nhãn dùng quy ước  $t \rightarrow t+1$  và căn chỉnh thời gian chặt chẽ.

### 3.4 Chỉ số đánh giá

RV: RMSE/MAE; Regime: ACC/F1-macro/nhị phân (Low vs Mid/High); Jump: AUROC/AP; Fade-classifier: ACC/F1/AUROC/AP; Regressor  $r^{\text{oc}}$ : RMSE/MAE/Corr. Backtest: hit-rate, Sharpe (252), P&L ròng sau phí.

### 3.5 Kết quả chính

**Model A (v2).** Tập huấn luyện/val/test lần lượt 968/58/61 ngày.

- **RV<sup>PK</sup>:** Val RMSE 0.0046 / MAE 0.0039; Test RMSE 0.0093 / MAE 0.0063.
- **Regime3:** Val ACC 0.81 / F1<sub>macro</sub> 0.45; Test ACC 0.28 / F1<sub>macro</sub> 0.15. Bản nhị phân (Low vs còn lại): Val 1.00; Test 0.72.
- **Jump:** Val AUROC 0.593 / AP 0.135; Test AUROC 0.702 / AP 0.378 (tỷ lệ dương  $\sim 6\text{--}10\%$ ).

**Model B (small-gap).** Tập huấn luyện/val/test lần lượt 417/21/23 ngày.

- **Classifier (fade):** Val ACC 0.67 / F1 0.74 / AP 0.56; Test ACC 0.57 / F1 0.62 / AUROC 0.56 / AP 0.62.

- **Regressor  $r^{\text{oc}}$ :** Val RMSE 0.0069 / MAE 0.0057 / Corr 0.18; Test RMSE 0.0073 / MAE 0.0050 / Corr 0.14.

### 3.6 Backtest quy tắc

Chúng tôi mô phỏng chính sách chỉ giao dịch *small-gap*, áp công chất lượng (gates) từ validation:  $p_{\text{fade}} \geq 0.66$ ,  $r_{\text{pred}}^{\text{oc}} \leq 0.0025$  (đối với fade),  $|\text{sent}| \leq 0.10$ ,  $z_{\text{news},w} \leq 1.8$ ,  $z_{\text{ent},w} \leq 0.8$ , chế độ *contrarian* với  $\text{fade30} \geq 0.52$ , phí 5bps. Kết quả tiêu biểu:

- **Opt A sweep (fade-only):** 15–17 lệnh trên  $\sim 536$  ngày, Sharpe  $\sim 0.04$ – $0.16$  (dương nhẹ), hit-rate  $\sim 0.53$ – $0.59$ .
- **Cấu hình hẹp (min\_retoc\_fade, min\_edge\_bps):** 24 lệnh/461 ngày, Sharpe  $\sim 8.0$ ; **cảnh báo:** mẫu nhỏ, biến động ước lượng cao  $\Rightarrow$  cần xác nhận tiền giao dịch/ngoại mẫu.

Phân tích theo bucket cho thấy hiệu quả chủ yếu nằm ở *small-gap*; *micro* và *large* có Sharpe âm khi xét toàn bộ.

### 3.7 Kiểm thử online

**Dịch vụ PhoBERT** (FastAPI, GPU): độ trễ sau warm-up quan sát  $\sim 24$ – $70$  ms/lệnh, p95 ban đầu  $\sim 1.9$  s do khởi động mô hình. **Bộ gom đặc trưng** ngày tạo `state/features_daily.parquet` theo thời gian thực. **Drift monitoring** (baseline 180d vs recent 60d): PSI cho `sent_score_mean`  $\approx 0.096$  (ổn), `sent_entropy_std_z_w`  $\approx 0.192$  (nhẹ), `news_count_z_w`  $\approx 1.43$  (cao) phản ánh khác biệt cỡ mẫu và nhịp tin gần đây; đã được cảnh báo trong dashboard.

### 3.8 Hạn chế & bài học

(i) Tín hiệu *follow* hiếm trong khung *small-gap*; (ii) Mô hình regime đa lớp suy giảm ngoài mẫu  $\Rightarrow$  sử dụng thêm bản nhị phân để định cỡ vị thế; (iii) Kết quả rất nhạy với cổng/gating và phí giao dịch; (iv) Cần baseline dài ( $\geq 180$ d) để ổn định  $z$  và giám sát drift.

## 4 Kết quả

### 4.1 Đánh giá ngoại mẫu (Phase 4)

Chúng tôi báo cáo kết quả trên tập **val** (2025-01-02  $\rightarrow$  2025-03-31) và **test** (2025-04-01  $\rightarrow$  2025-06-30), chi phí giao dịch được xem xét ở phần backtest.



#### Model A (Risk gauge, v2).

- **Realized Volatility (Parkinson,  $RV^{PK}$ ):** Val RMSE 0.0046 / MAE 0.0039; Test RMSE 0.0093 / MAE 0.0063.
- **Regime 3 lớp:** Val ACC 0.81 /  $F1_{macro}$  0.45; Test ACC 0.28 /  $F1_{macro}$  0.15. Biểu thể *nhị phân* (Low vs Mid/High): Val 1.00; Test 0.72.
- **Jump (hiếm,  $\sim 6\text{--}10\%$ ):** Val AUROC 0.593 / AP 0.135; Test AUROC 0.702 / AP 0.378.

#### Model B (Small-gap decision).

- **Classifier (fade/không):** Val ACC 0.67 / F1 0.74 / AUROC 0.48 / AP 0.56; Test ACC 0.57 / F1 0.62 / AUROC 0.56 / AP 0.62.
- **Regressor  $r^{oc}$ :** Val RMSE 0.0069 / MAE 0.0057 / Corr 0.18; Test RMSE 0.0073 / MAE 0.0050 / Corr 0.14.

## 4.2 Backtest quy tắc (Phase 5)

Chiến lược chỉ giao dịch *small-gap* với các cổng chất lượng (gates) rút ra từ validation:  $p_{fade} \geq 0.66$ ,  $\min\_retoc\_follow = 999$  (bỏ follow),  $\max\_retoc\_fade \in [0.0020, 0.0025]$ ,  $|sent| \leq 0.10$ ,  $z_{news,w} \leq 1.8$ ,  $z_{ent,w} \leq 0.8$ , chế độ *contrarian* với  $fade30 \geq 0.52$  (cửa sổ 30 ngày). Phí giao dịch mặc định 5bps.

- **Opt A (sweep):** *fade-only*,  $n=15 \sim 17$  lệnh trên  $\sim 536$  ngày, hit-rate  $\sim 0.53\text{--}0.59$ , Sharpe  $\sim 0.04\text{--}0.16$  (dương nhẹ, bền hơn qua tham số).
- **Opt A (hẹp hơn)** với  $\min\_retoc\_fade$  và  $\min\_edge\_bps$ :  $n=24$  lệnh/461 ngày, Sharpe  $\sim 8.0$ ; **lưu ý:** cỡ mẫu nhỏ, cần xác nhận ngoại mẫu.
- **Opt B (follow-only)** cho *small-gap* cho kết quả âm (Sharpe  $\ll 0$ ), phù hợp quan sát thiên hướng *mean reversion* ở *small-gap*.

Phân tích theo bucket cho thấy lợi nhuận tập trung ở *small-gap*; *micro* và *large* cho Sharpe âm khi xét rộng.

## 4.3 Phục vụ trực tuyến & giám sát (Phase 6)

Dịch vụ PhoBERT (FastAPI, GPU) trả lời sau *warm-up* với độ trễ  $\sim 24\text{--}70$  ms/lệnh (p95 khởi động  $\sim 1.9$  s). Bộ gom *online* tạo *state/features\_daily.parquet* theo ngày từ log. Ví dụ *live\_risk.json* ngày 2025-09-02: *regime3\_pred* = 1 (mid), *size\_today* = 0.7, *jump\_prob*  $\approx 1.8 \times 10^{-5}$ . Với *gap* = 0.0043 cùng bộ cổng “opt A sweep”, *live\_signal.json* ra *decision* = *skip* do vi phạm các cổng:  $|sent| > 0.1$ ,  $p_{fade} < 0.66$ ,  $r_{pred}^{oc} > 0.0025$ .

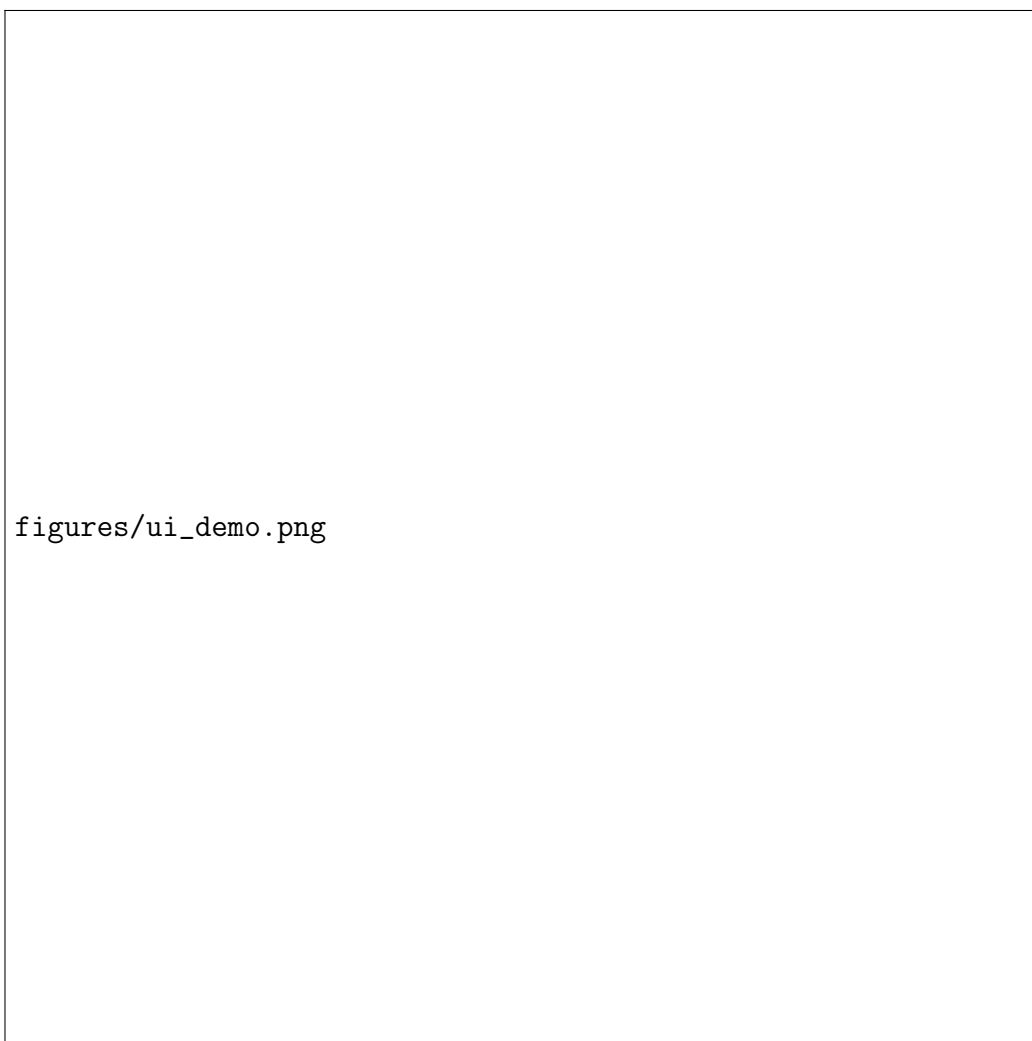
**Drift monitoring (180d vs 60d).** PSI cho `sent_score_mean`  $\approx 0.096$  (ổn), `sent_entropy_std_z_w`  $\approx 0.192$  (nhẹ), `news_count_z_w`  $\approx 1.43$  (cao do khác biệt cỡ mẫu/nhịp tin). Dashboard sinh `metrics.json` và `report.html` để theo dõi hàng ngày.

#### 4.4 Minh hoạ giao diện người dùng

Hình 1 minh hoạ giao diện demo: nhập tin, xem phân tích sentiment, risk gauge (regime/jump) và tín hiệu giao dịch cùng vết cổng (gate audit).

#### 4.5 Tổng kết

(1) Tín hiệu *fade* ở *small-gap* mang lại tỉ lệ thắng ổn định khi áp cổng chất lượng; (2) Mô hình risk (A) hữu ích để *size* vị thế và cảnh báo ngày rủi ro/jump; (3) Chuỗi phục vụ & giám sát realtime vận hành tốt, đã có cảnh báo drift; (4) Cần mở rộng ngoại mẫu, tăng số lệnh đủ lớn và chuẩn hoá phí thực tế trước khi triển khai giao dịch.



figures/ui\_demo.png

Hình 1: Giao diện demo: luồng  $News \rightarrow Sentiment/Entropy \rightarrow Risk\ gauge \rightarrow Trading\ signal$  với gate audit.