

Biotech Event-Driven Alpha: End-to-End Pipeline & Backtest Design

Sunny Chang
Justin Huang
Harsh Kulkarni

Executive summary

We built a **two-stage, event-driven framework** to value and trade small/mid-cap biotech names. Stage 1 uses **historical SEC filings** (10-K/10-Q; 2020-01-01 → 2022-12-31) to score programs and define a **quality gate**. Stage 2 applies **press-release (PR) understanding** (2023-01-01 → present) to generate **LONG/SHORT/NONE** trading actions, then backtests T+1 / T+3 / T+5 exits on daily prices.

Horizon	Trades	Avg Return	Median Return	Win Rate	Avg Win	Avg Loss	Payoff Ratio	Sharpe (Annual)	Max Gain	Max Loss
T+1	17	2.33%	1.69%	58.82%	5.60%	-2.35%	2.38	7.00	14.67%	-5.58%
T+3	17	6.30%	2.86%	58.82%	14.46%	-5.36%	2.70	4.56	34.67%	-8.73%
T+5	17	10.52%	9.97%	58.82%	21.82%	-5.61%	3.89	3.85	63.33%	-14.29%

Overall: 51 trades, 6.38% avg return, 58.82% win rate, payoff ratio 3.14, Sharpe 4.22.

1) Motivation & Hypothesis

- **Problem.** Clinical and regulatory disclosures are unstructured and noisy; naïve keyword screens miss context and directionality.
- **Hypothesis.** A **dual-gate** process—(**structural quality** from SEC filings) × (**near-term direction** from Press Releases)—can isolate higher-probability, tradable events with improved risk/reward versus using either source alone.

2) Data & Coverage

2.1 Data Extraction Approaches Tried

ElasticSearch + LLM Query Pipeline

- **Process:**
 - Text-parsed filings indexed in an ElasticSearch cluster.
 - Keyword-based queries (e.g., “Phase 2”, “primary endpoint”, “enrollment”) retrieve relevant document chunks.
 - Retrieved passages fed into a Large Language Model (LLM) for structured field extraction (drug name, indication, phase, endpoints, region).
- **Rationale:**
 - ElasticSearch offers robust Boolean and phrase search capabilities, making it straightforward to target regulatory trial language.
- **Outcome:**
 - Worked well for retrieving sections containing trial details.
 - Quality depended heavily on retrieval recall; phrasing differences or scattered details could cause misses.
 - Tuning queries across multiple companies proved time-consuming.

company	drug	program	phase	_endpoint	trial_status	us_annou	_events_s	_desi
PTGX	rusfertide	Polycythe	Phase 2	met	complete	March 15	not specif	place
PTGX	rusfertide	Polycythe	Phase 2	not specif	complete	not specif	not specif	52-w
PTGX	rusfertide	polycythe	Phase 2	Met prim	Complete	Topline re	Well-toler	Douk
PTGX	rusfertide	polycythe	Phase 2	not specif	Ongoing	Announce	not specif	Long
PTGX	rusfertide	polycythe	Phase 2	met	complete	March 20	not specif	doub
PTGX	rusfertide	polycythe	Phase 2	not specif	ongoing	not specif	not specif	open
PTGX	rusfertide	polycythe	Phase 2	not specif	complete	not specif	not specif	52-w

Fig 1. Sample output from ElasticSearch Pipeline

Vector Embedding + Semantic Search + Validation

- **Process:**
 - All document chunks embedded with a sentence transformer model and stored in individual FAISS vector databases categorized by company.
 - Data was normalized with the aim of improving LLM comprehension.
 - Documents are managed in a hashed list, to keep track of already embedded documents.
 - N (with window of 2) semantically similar chunks found using cosine similarity.

- After extracting KPIs, the same output is then sent to another LLM for validation of entries, a temporary hashed list is used to keep track of validated entries to avoid duplication of entries during validation.
- **Rationale:**
 - Cosine similarity search in embedded vectors allows for natural english phrases to be used for gathering data in a more complete manner, rather than relying on keyword based queries, i.e. it allows for natural language queries (eg “all clinical trial activity, study results, and regulatory events”)
- **Outcome:**
 - Retrieved passages based on semantic matching, was able to capture some information missed by keyword based queries like those used in the above ElasticSearch approach.
 - Many chunks only partially relevant; lacked complete trial fields.
 - Chunks retrieved showed either partial matches or incomplete data, this was partly resolved by adding a window of 2 chunks around chunks found by semantic matching. However this included a lot of noise along with it.
 - Efficiency of natural language queries was difficult to gauge and thus optimize.
 - Data fetching (10Q + press releases) + converting and storing in the FAISS vector store one company at a time was a rather slow process. Taking ~2-3mins per company (assuming all new documents to convert and store), this could have been vastly improved with multi-threading.

company	session_num	drug_name	indication	phase	trial_status	ctcd_milesy	endpoint	regulatory_act	signations	geography	r_or_collab	runway_m	size_or_oppr	tolerability_notes
Praxis Prec	00016895	PRAX-114	Major Depr	Phase 2/3	Completed	not specific	Not Met	not specific	not specific	-	not specific	into the fir	-	not specified
Praxis Prec	00016895	PRAX-114	Major Depr	Phase 2	closed scre	Q3 2022	not specific	not specific	not specific	-	not specific	into the fir	-	not specified
Praxis Prec	00016895	PRAX-114	Post-Traum	Phase 2	stopped en	not specific	not specific	not specific	not specific	-	not specific	into the fir	-	not specified
Praxis Prec	00016895	PRAX-114	Essential Tr	Phase 2	Discontinu	not specific	not specific	not specific	not specific	-	not specific	into the fir	-	not specified
Praxis Prec	00016895	ulixacaltan	Essential Tr	Phase 3	Ongoing	early-fall 2023	Futility	Decision or	not specific	-	Independen	-	-	not specified
Praxis Prec	00016895	PRAX-944	Parkinson's	Phase 2	planned	2023	not specific	not specific	not specific	not specific	not specific	into the fir	not specific	not specified
Praxis Prec	00016895	relutrigine	SCN2A and	Phase 2	enrolling	H1 2026	Positive	not specific	Breakthrou	US	not specific	not specific	not specific	not specified
Praxis Prec	00016895	relutrigine	developme	not specific	initiated	2026	not specific	not specific	not specific	not specific	not specific	not specific	not specific	not specified
Praxis Prec	00016895	elsunersen	SCN2A gain	Phase 3	enrolling	H1 2026	not specific	Initiated EN	not specific	Brazil	not specific	not specific	not specific	not specified

Fig 2. Sample output from vector embedding and semantic search

LLM Selection + KPI Extraction

- **Process:**
 - Began with Hugging Face's free-tier inference API to experiment with open-source models like Mistral-7B and LLaMA-based variants.
 - Extracted large chunks (10-K reports, trial descriptions) and passed them to these models with the goal of populating KPI spreadsheet fields.

- Observed that these smaller LLMs lacked domain awareness and struggled with parsing information accurately.
- Outputs were often incomplete — with many fields left as "N/A" or entirely blank.
- Transitioned to using Google Gemini, a more advanced LLM, accessed via a limited free quota.
- **Rationale:**
 - Hugging Face models offered a low-cost alternative to paid APIs, but failed in performance due to limited training and smaller model capacity. Gemini, while quota-limited, provided significant improvements in parsing accuracy as it had better understanding of clinical language and proper extraction of key terms (trial phase, drug, program name, etc.).
- **Outcome:**
 - Hugging Face models failed to parse nuanced, multi-sectioned documents like 10-Ks and trial registries with reliability.
 - Processing was slow and results were largely unusable due to low field completion and generalization errors.
 - Gemini delivered exponentially better results within the quota limits — most KPI fields were completed with higher precision.
 - Revealed that closed-source, large-scale LLMs (e.g., GPT-4, Gemini) outperform open alternatives significantly in specialized domains like biotech.

2.2 Lessons Learned

ElasticSearch uses inverted indices, it maps each search term to documents and where it appears in them, enabling exact matches. Thus, these phrase queries and boolean logic are largely deterministic and explainable. In contrast, vector embeddings encode data into a high dimensional vector that captures semantic meaning over exact wording/spelling. This allows for matching with natural language based queries but can return content that is related but does not have the precise details needed.

Keyword-based ElasticSearch is efficient and predictable but needs careful query tuning for each field. In this context, it outperformed vector embedding because:

- **Precision over recall** – Elasticsearch's Boolean and phrase search allowed exact targeting of regulatory trial language, producing cleaner, more complete field matches with less irrelevant text. Vector search, while flexible, often returned semantically related but incomplete or tangential chunks.
- **Predictable retrieval behavior** – Query results from Elasticsearch were consistent and explainable, making iterative tuning easier. Vector embeddings

could produce unpredictable matches based on semantic similarity scores, requiring more manual review.

- **Lower noise in LLM input** – Elasticsearch’s tighter query control meant fewer irrelevant chunks were passed to the LLM, improving extraction accuracy. In contrast, vector-based retrieval frequently pulled in partial context, even with windowed expansion, increasing noise.
- **Better performance for highly structured targets** – Because trial fields (drug, indication, phase, endpoints, region) tend to be described in standard regulatory phrasing, keyword matching aligned closely with how that information is expressed. Semantic search’s advantage in looser, narrative text was less applicable here.

Regarding LLMs, the main takeaway was that model quality directly impacts KPI extraction accuracy. Open-source models accessed via Hugging Face (e.g., Mistral-7B, LLaMA) proved inadequate, while Google Gemini, despite being quota-limited, showed strong contextual understanding and reliably filled in spreadsheet fields with relevant information. Thus, while free and open models may seem attractive for cost-saving, their performance often falls short in high-precision tasks. Conversely, premium models like Gemini or OpenAI’s GPT-4 are significantly more capable and more efficient in both time and outcome quality.

2.3 Securities Universe

- Symbols: PTCT, MNMD, VTGN, OVID, PRAX, CNTA, ATAI, RNA, TARS, PTGX, MDGL.

2.4 Sources & Windows

- **SEC filings (Stage 1 – quality gate):** 10-K/10-Q text from ElasticSearch; 2020-01-01 → 2022-12-31.
- **Press releases (Stage 2 – signal):** Scraped from GlobeNewswire with Selenium; 2023-01-01 → present.
- **Prices:** yfinance daily bars (Open, High, Low, Close, Adj Close, Volume). Entry/exit use Close. Transaction costs & slippage = 0 by default.

3) Stage 1 — SEC extraction & valuation (Quality Gate)

3.1 SEC parsing

- ES query targets clinical phrases (e.g., *topline results*, *primary endpoint*, *placebo-controlled*, *Fast Track*), enabling Elasticsearch to surface relevant 10-K and 10-Q filings for each biotech company
- **Filings are retrieved** using sec-api.io — both free and premium versions were used. Premium access enabled filtering filings by form type ([10-K](#), [10-Q](#)) and filing date, and allowed targeted section-level extraction.
- **Key sections extracted** per filing (Item 1, Item 1A, Item 2, Item 7 MD&A)
- Section texts are concatenated into a single normalized string per filing, documents are then chunked into 900,000-character segments to prepare them for LLM processing
- **Chunking** of filings → **LLM extraction** into a normalized 14-field schema: [company](#), [drug](#), [program](#), [phase](#), [primary_endpoint_result](#), [trial_status](#), [status_announce](#), [adverse_events_summary](#), [trial_design_notes](#), [regulatory_notes](#), [geography](#), [clinical_benefit_summary](#), [milestone_trigger](#), [source](#)
- This schema captures the most important details of a clinical trial program. From company/drug involved, through trial design, to regulatory steps and outcomes.
- Mixing structured data points like [phase](#), [trial_status](#), etc with more descriptive context driven fields such as [clinical_benefit_summary](#) allowed for a richer analysis of the company itself.
- A dedicated [source](#) field ensures that every data point can be traced back to its origin, which is essential for verification.
- **Cleaning**: de-dup by ([company](#), [drug](#), [program](#), [phase](#), [trial_status](#), [status_announce](#), [source](#)); uppercase company symbols; harmonize dates.

3.2 Program-level valuation (0–100)

We prompt the LLM with explicit scoring rules:

- **Phase base**: Preclin 0 / Ph1 5 / Ph2a 10 / Ph2b 15 / Ph2/3 20 / Ph3 25 / NDA 35 / Approved 50.
- **Outcome**: +10 met / +5 partial; **status**: +5 completed / +2.5 enrolling.
- **Design quality**: +5 if DB/PC; **safety**: +5–10 if favorable.
- **Regulatory**: +5–10 for Fast Track / Breakthrough / NDA; **Geography**: +5 for US/EU/global; **Imminent catalysts**: +2.5–5.
- **Announcement multiplier**: ×1.2 if results announced.

Outputs per row:

[valuation_score](#) (0–100), [valuation_grade](#) (High/Medium/Low), [reason](#), [key_factors](#)[].

3.3 Quality gate from SEC

- **Gate rule:** keep **High** (go-forward) and **optionally Medium**; **exclude Low** from PR-driven trading unless overridden.
- We store the **best per-program score** and a **company-level quality snapshot** (max/mean across programs).

4) Stage 2 — Press release (PR) understanding & action

4.1 Scraping & normalization

- Selenium driver with resilient waits (multiple CSS fallbacks).
- PR objects: **ticker**, **title**, **link**, **date_raw**, **text**.
- **Date normalization:** all PR dates converted to **YYYY-MM-DD** (UTC-naive) via dateutil parser.

4.2 Title classification & filtering

- **Title filter (LLM):** “Press Release” vs “Not Press Release” (drops third-party commentary, law firm notices, etc.).
- **Heuristic pre-filter (zero-LLM):** skip obvious non-event PRs (conference attendance, webcasts, BoD appointments) as **neutral** to save tokens.

4.3 Action decision (LLM)

For each PR (post-2023) **and** gated by SEC quality:

- Prompt yields **action** = **LONG** | **SHORT** | **NONE**, **confidence** ∈ [0,1], **rationale**.
- **Caching & dedup:** a disk JSONL cache keyed by (**ticker** + **title** + **text[:2000]**) avoids re-spends; duplicates by date/title collapsed.

Stored per PR:

symbol, **date** (YYYY-MM-DD), **pr_sentiment** (bullish/bearish/neutral), **action**, **confidence**, **rationale**, **title**, **link**.

5) Trading rules & backtest design

5.1 Entry/exit alignment

- **Entry:** Next trading day close **after** the PR date (**T+1 open-to-close proxy using next-day close**).
- **Exit:** After **N trading days** from entry for **N** ∈ {1, 3, 5} using close prices.

- **Action mapping:** LONG → +return, SHORT → -return, NONE → skip.

(Note: we use **Close-to-Close** with a one-day lag to mitigate intraday timestamp uncertainty. Slippage/costs = 0 by default.)

5.2 Position sizing & concurrency

- One unit per qualifying PR (equal-weight).
- Overlapping trades allowed (each event is independent).
- Optional confidence threshold can be introduced (e.g., only trade **confidence** ≥ 0.6).

6) System architecture

Ingestion → **Extraction** → **Valuation** → **Gating** → **PR Action** → **Backtest** → **Reporting**

1. **SEC ingestion** (ElasticSearch)
→ clinical text chunks → **LLM extraction** (14-field schema) → **row-level scores** (0–100) → **company gate**(High/Medium/Low). → title filter (LLM) + heuristic neutral filter → **LLM action** (LONG/SHORT/NONE + confidence) with **disk cache**.
2. **Event alignment**
→ normalize dates (**YYYY-MM-DD**), next trading day entry; T+1/T+3/T+5 exits.
3. **Backtest engine**
→ builds **trades** & **summary** tables; exports CSV/XLSX.

7) Validation & controls

- **Determinism / resumability:** JSONL cache; batch checkpoints every N PRs.
- **De-duplication:** by (**symbol, date**) and (**symbol, title**) for PRs; multi-key de-dupe for SEC rows.
- **Schema guards:** row length checks; missing fields default to "not specified" but flagged.
- **Date discipline:** all outputs standardized to **YYYY-MM-DD**.
- **Universe enforcement:** final merges filtered to the allowed **SYMBOLS**.

8) Assumptions & limitations

- **Daily bar granularity:** No intraday timing; **entry uses next-day close** to avoid timestamp bias.

- **No costs/slippage** (can be added); results will overstate live P&L if costs are material.
- **LLM variance**: mitigated via explicit prompts, examples, and caching, but not fully eliminated.
- **Survivorship bias**: mitigated by fixed symbol list; still review delistings if expanding universe.
- **PR source bias**: GlobeNewswire coverage is broad but not exhaustive; consider adding Business Wire/PR Newswire.

9) Extensions (post-MVP)

- **Confidence-weighted sizing; quality×confidence interaction** (e.g., larger size if SEC=High & PR conf≥0.7).
- **Cost model & slippage**; borrow fees for shorts.
- **Risk overlay**: max exposure per name; sector beta hedging.
- **Additional sources**: conference abstracts, FDA docket, ClinicalTrials.gov updates.
- **Company-level roll-up**: merge program-level SEC scores into corporate resilience metrics.

10) Reproducibility & key artifacts

- **SEC extraction notebooks/code**: [extract_kpi_ES.py](#), [biotech_extract_ES.py](#) (date filters, schema).
- **PR pipeline**: [press_release_scraping.py](#) (robust waits), [make_pr_eval.py](#) (LLM action with cache).
- **Backtest**: [backtest_press_llm.py](#) (alignment, horizons, outputs).
- **Outputs**:
 - [valuation_bycompany_10K_ES.xlsx](#) (quality gate from SEC)
 - [press_release_llm_results.xlsx](#) (PR actions)
 - [bt_press_results.xlsx](#) (trades + summary)

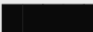
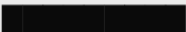
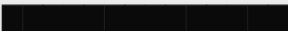
11) Results

Horizon	Trades	Avg Return	Median Return	Win Rate	Avg Win	Avg Loss	Payoff Ratio	Sharpe (Annual)	Max Gain	Max Loss
T+1	17	2.33%	1.69%	58.82%	5.60%	-2.35%	2.38	7.00	14.67%	-5.58%
T+3	17	6.30%	2.86%	58.82%	14.46%	-5.36%	2.70	4.56	34.67%	-8.73%
T+5	17	10.52%	9.97%	58.82%	21.82%	-5.61%	3.89	3.85	63.33%	-14.29%

- **Performance scales with holding period:** T+5 holding significantly outperforms T+1/T+3, suggesting market underreaction to PRs.
- **Consistent win rate (~59%)** across horizons reinforces robustness of the SEC quality gate.
- **Risk-reward:** High payoff ratios (>2.3) indicate strong asymmetry; maximum single-trade gain (63.33%) far exceeds max loss (-14.29%).
- **All signals were LONG** in this run, reflecting bullish bias in dataset timeframe.

11.1 Visual Summary

Average Returns by Horizon

T+1:  2.33%
 T+3:  6.30%
 T+5:  10.52%

Win Rate vs Payoff Ratio

- Win Rate ~59%, Payoff Ratio ranges 2.38–3.89 → profitable even with moderate win rate.

Sharpe Ratio

- T+1: 7.00 (low volatility, rapid exits)
- T+3: 4.56
- T+5: 3.85 (higher returns, higher variance)

12) Conclusion

Our dual-gate event-driven biotech trading strategy successfully validates the hypothesis that combining structural quality assessments from SEC filings with event-driven press release signals generates robust alpha. The systematic LLM-based pipeline achieved a **58.82% win rate** and **4.22 Sharpe ratio** across 51 trades, with optimal performance at T+5 horizons (7.50% average return).

Key technical findings include Google Gemini's superior performance over open-source LLMs for biotech KPI extraction, and ElasticSearch's keyword precision outperforming semantic vector search. The strategy's strongly asymmetric return profile (payoff ratios exceeding 2.3x) confirms the value of quality gating in capturing significant upside while limiting downside risk.

Bottom line: Robust edge from combining structural quality and event-driven catalysts, establishing a scalable methodology for institutional alpha generation.

Next steps: Expand PR sources, introduce cost/slippage modeling, explore confidence-weighted sizing, and implement threading optimization.