



Martin Hynie
LaunchPT.com



RETURN OF THE EXPLORER

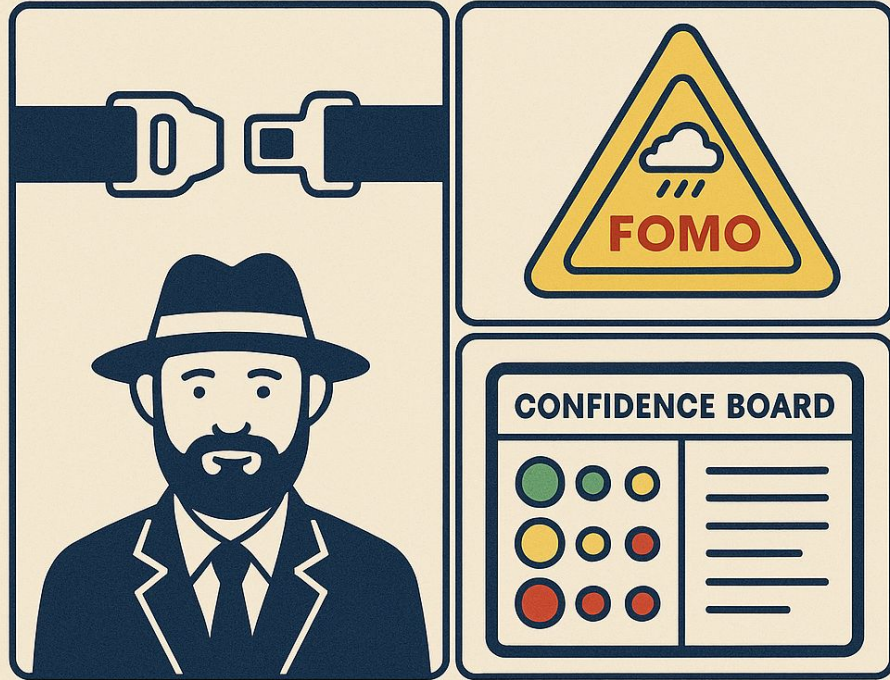
HITL

KILL SWITCH

IN-FLIGHT SAFETY FOR AI SYSTEMS



PLEASE FASTEN YOUR GUARDRAILS



CONSEQUENCES > VIBES

Confidence isn't a feeling;
it's evidence we can show."



FOMO GRAVITY



**HYPE picks directions;
evidence sets speed.**



Build lenses... Not laws or limits.
Use them to notice “more”.

ACT I: WHY THINGS FAIL NOW



INTERACTION

**DO NOT
INCORRECTLY
INTERACT**



INFRASTRUCTURE

**BROKEN SYSTEMS
MAY OCCUR**



INSTITUTION

**AVOID
OVERWHELMING
PRESSURE**

AIRPORT

29.01.11

★ USA ★

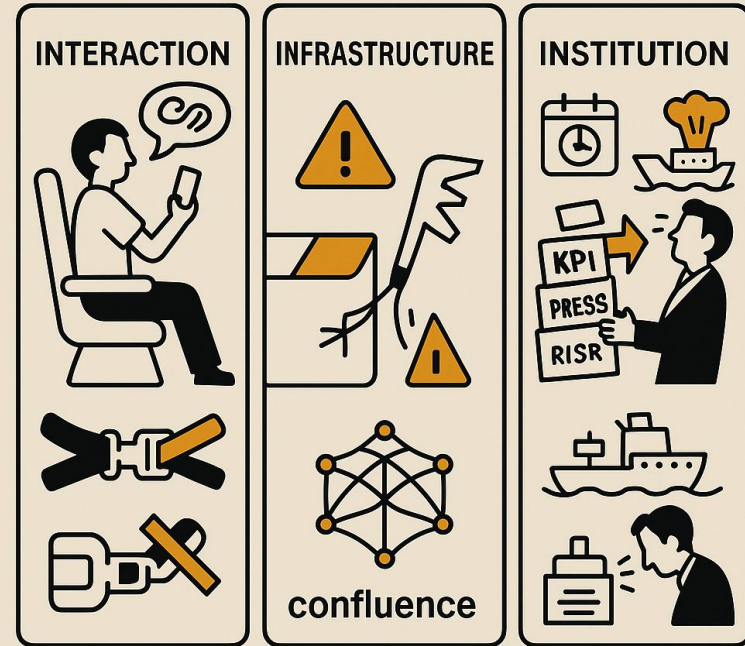
Interaction: prompt/context drift, retrieval variance, misalignment.

Infrastructure: fuzzy data lineage, stale evals (constant risk), degraded guardrails, missing traces.

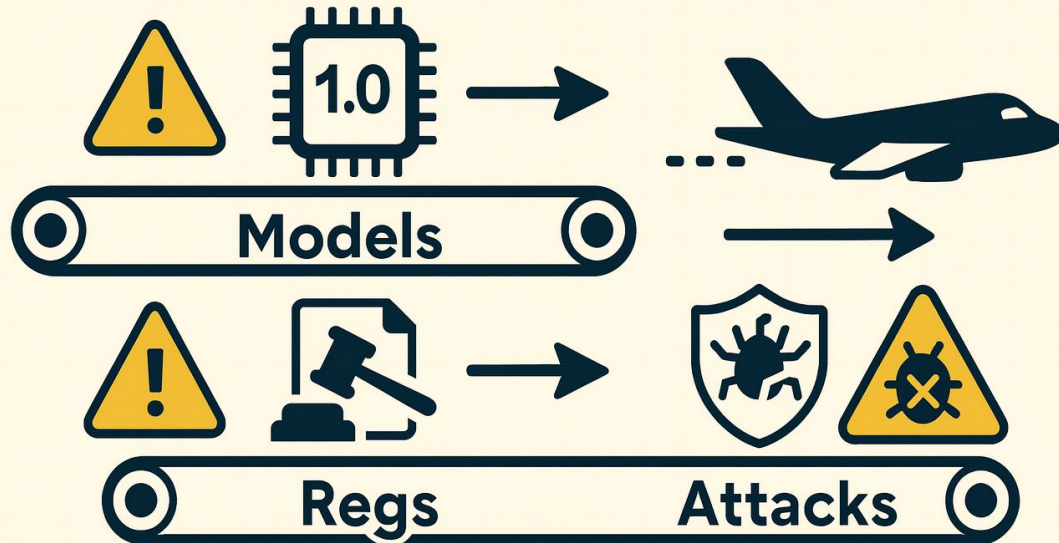
Institution: incentives, process debt, systemic shipping pressure, fomo

Complexity 101: outcomes emerge from interactions, not single rca

THREE LENSES



NON-STATIONARY WORLD

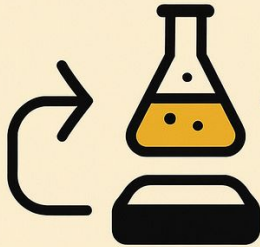


Fixing yesterday's cause \neq shrinking tomorrow's risk surface.

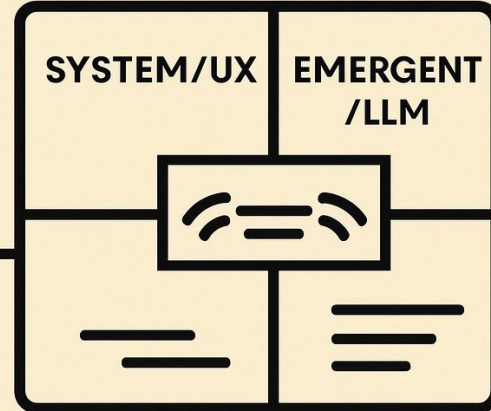


ARTIC
OCEAN

ACT II—BUILD TO EXPLORE (AND TEST ON TWO PLANES)



QE BRIDGE



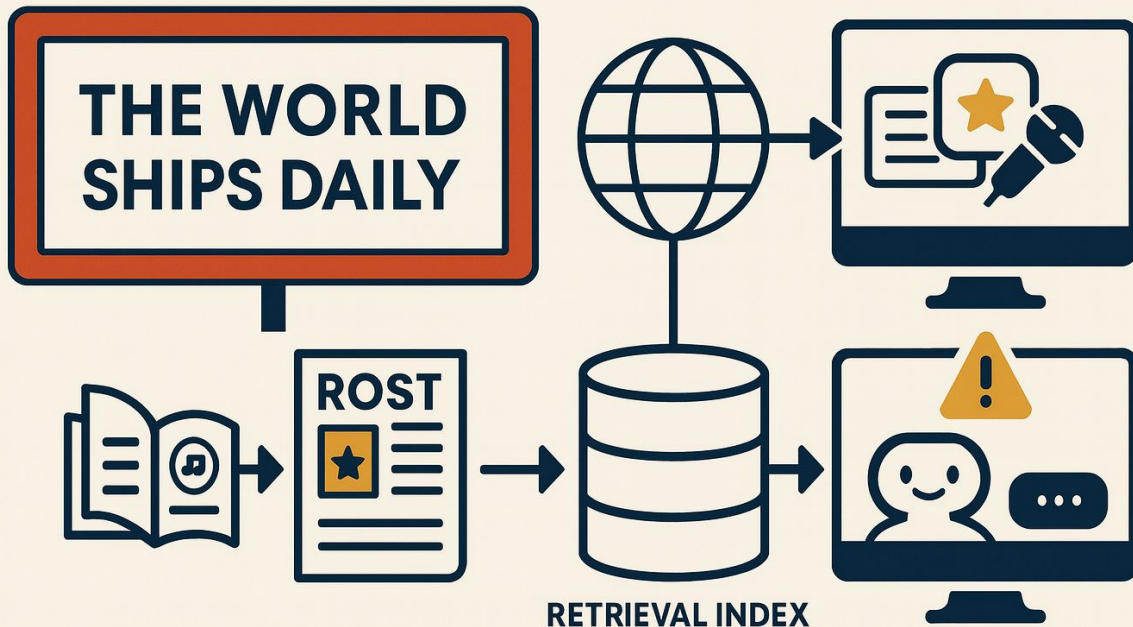
PACIFIC
OCEAN

OCEANIA

AUTRALIA



REALITY DRIFT



"YOUR MODEL DIDN'T CHANGE. THE WORLD DID."





LEARNING REVIEWS > BLAME

“If an incident doesn’t create new probes, we paid tuition and learned nothing.”

CONTRIBUTING-FACTOR TIMELINE



Infrastructure Interaction

CONTROLS CATALOG

 PRESENT

 ABSENT

 DEGRADED

OUTPUTS



NEW PROBES



TELEMETRY ADDS



TELEMETRAINTY BUDGT UPDATE

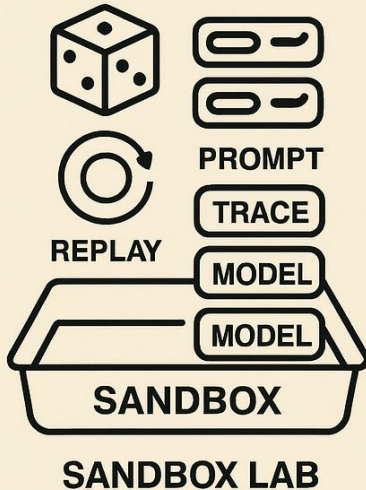




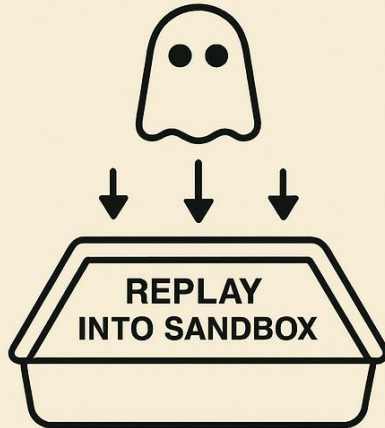
BUILD TO EXPLORE

Don't just test the system—build little systems to explore it.

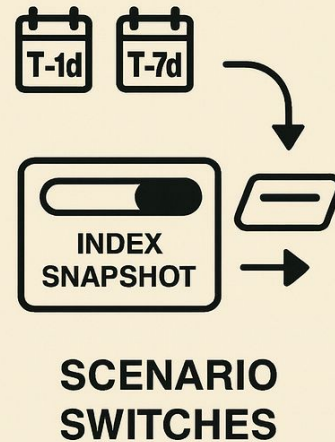
SANDBOX LAB



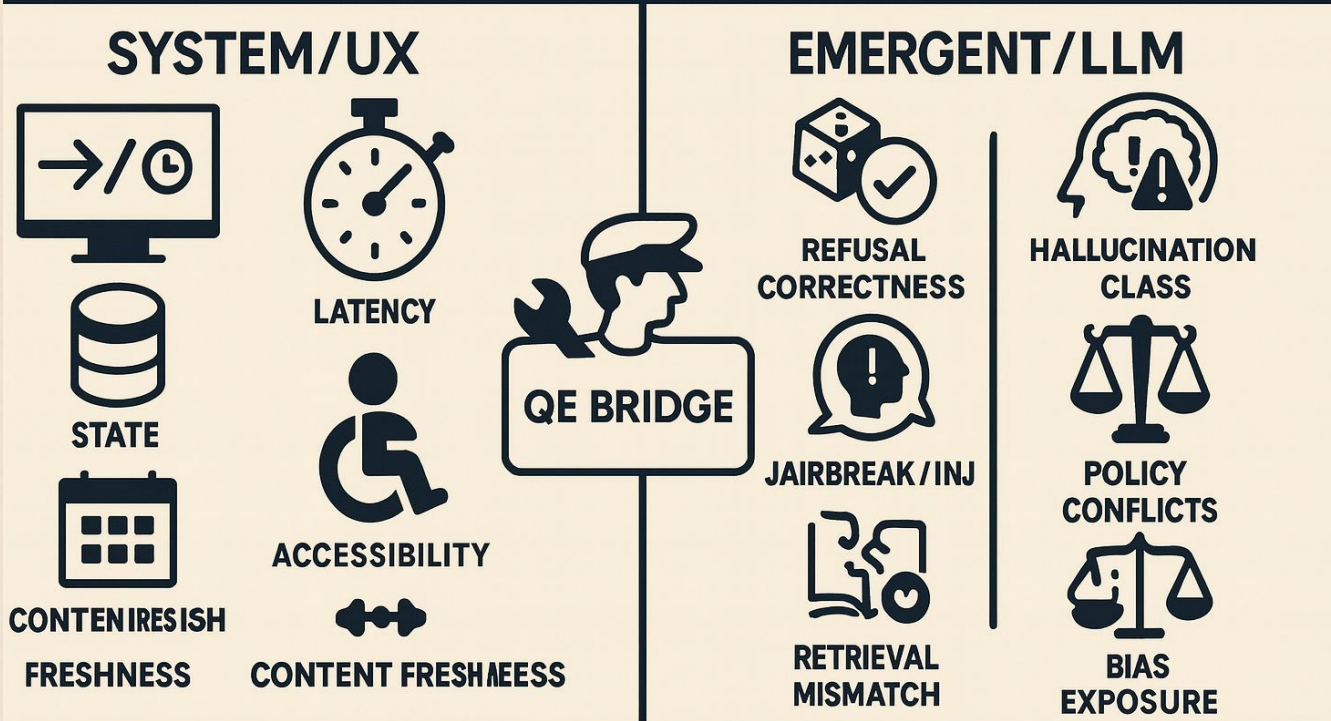
SHADOW TRAFFIC



SCENARIO SWITCHES

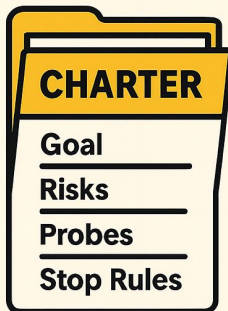


TWO TESTING PLANES

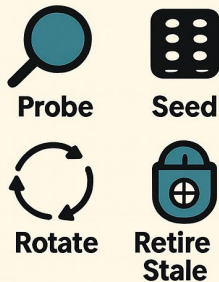
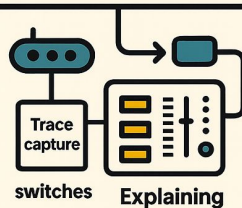




**Test cases check claims.
Charters chase unknowns.**



Uncertainty



**UNCERTAINTY
CLASSES**



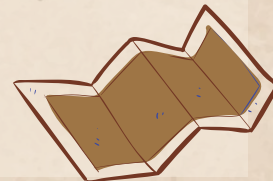
**Build systems that
explain themselves.**

As a tester, you have to advocate for building the system in such a way that these charters are explorable by design.

YOU...

have to be part of the design process.

No downstream tool will save you if you did not build a system expected to explain itself.





Reference + adversarial sets weighted by impact.

Track confidence deltas over time

Use data blast-radius constraints... so one trending source can't flip everything

GUARDRAILS + HITL & DRIFT CONTROLS



TIME-SLICED EVALS



RETRIEVAL VOLATILITY



CHANGEFEED
→ PROBES



QUARANTINE
/ SANDBOX

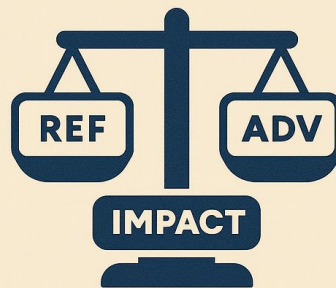


SNAPSHOT & FREEZE

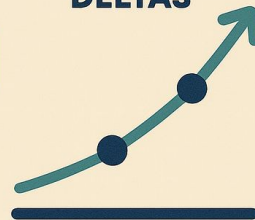
TIME-SLICED EVALS

RETRIEVAL VOLATILITY

HARNESS = MIX BY IMPACT



CONFIDENCE
DELTA



↓ Toxicity false-negatives

↑ Refusal correctness

Retrieval mismatch

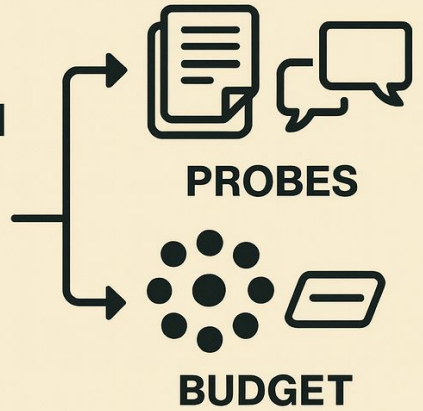


COVERAGE VS CONSEQUENCES



ADJACENT FIELDS VIGILANCE

SUFFIX ATTACKS

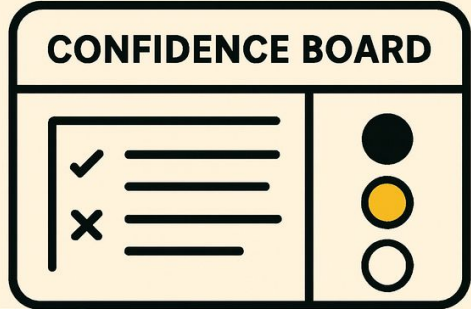
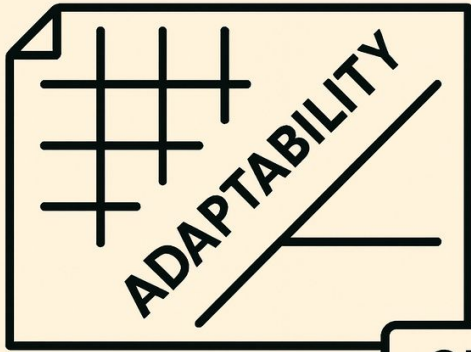


You're not security, but you are early warning.

"You don't have to know everything—just what changes your r risk today."

ARTIC
OCEAN

ACT III – QE = ARCHITECTURE & EXPLAINABILITY + GENEVAL LEADS DESIGN



QE BRIDGE

PACIFIC
OCEAN

OCEANIA

AUTRALIA



QE = ARCHITECTURE & EXPLAINABILITY

shape the system



Routing layers &
kill switches

ADAPTABILITY



Policy-as-code



Prompt/trace
capture



Uncertainty hooks



Explainability
hooks



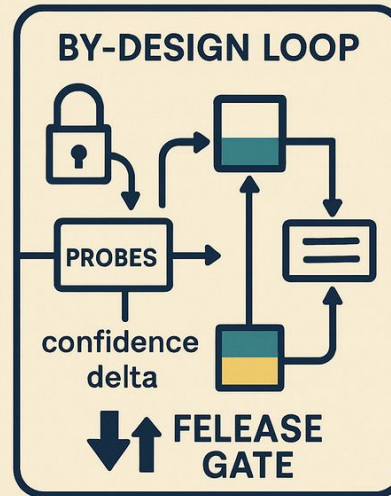
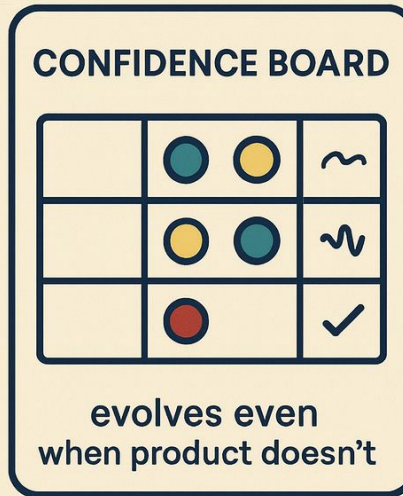
Explainability
hooks

We can't shrink the unknowns—but we can shape the system that meets them.



GenEval (By Design)

“Evaluation leads design.





CASE VIGNETTE:

INCIDENT → PROBES → GUARDRAILS → BUDGETS → TEEMETRY

INCIDENT



UNACCEPTABLE
OUTPUT

PROBES



REF + ADV
PARAPHRASES

GUARDRAILS & ROUTING



NEAR-ZERO

YOUTH/
DISTRESS

BUDGETS



TRACE



POLICY
APPLIED



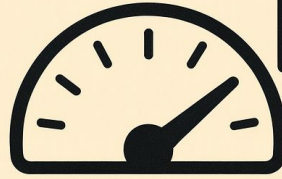
CONFIDENCE

TELEMETRY

MODEL HARDENS? TIGHTEN PROBES... WORLD SHIFTS? TIME-SLICE
EVALS, VOLATILITY SLO, QUARANTINE, SNAPSHOT/FREEZE.



**MAKE CONFIDENCE
VISIBLE**



MAKE CONFIDENCE VISIBLE

Q & A



Martin Hynie | LaunchPT.com