

*ICSE 2027 Paper Progress*

# Behavior Vector Coverage

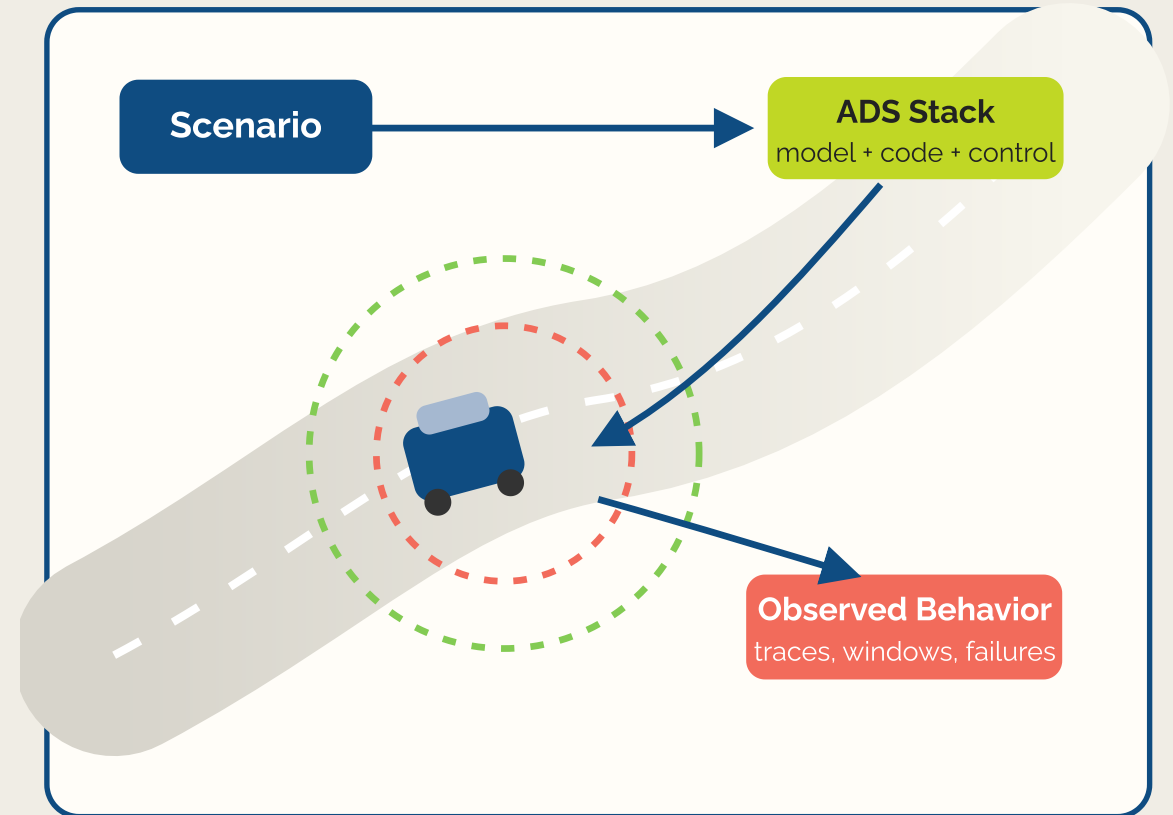
BVS Coverage Project

ADS test adequacy progress update from current implementation and InterFuser RQ1 artifacts

June 5, 2026

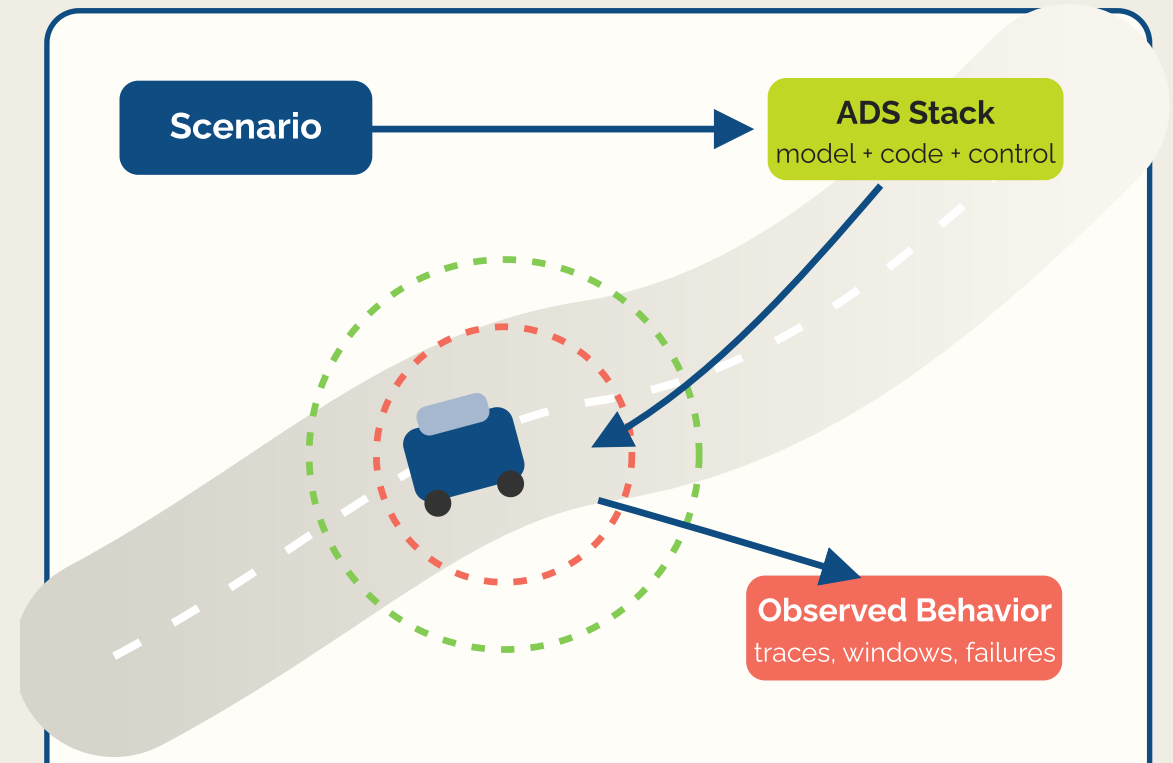
# ML-Enabled Testing Is Hard

- **ADS behavior emerges from model + code + traffic**
- **Failures are rare, contextual, and expensive to reproduce**
- **DNN coverage alone misses system-level behavior**
- **A scenario can pass while still exercising little behavior**



# ML-Enabled Testing Is Hard

- ADS behavior emerges from model + code + traffic
- Failures are rare, contextual, and expensive to reproduce
- DNN coverage alone misses system-level behavior
- A scenario can pass while still exercising little behavior



Adequacy should talk about exercised behavior, not only inputs.

# Scenario Inputs Drive ADS Testing

- **Simulation turns scenarios into executable inputs**
- **Search mutates maps, actors, routes, timing, weather**
- **Diversity is often defined over scenario parameters**
- **The goal is to expose more critical behaviors**

### Map

Town / lanes / junctions

### Actors

vehicles / pedestrians

### Timing

spawn / speed / route

### Weather

rain / fog / lighting

### Oracle

collision / lane / red light

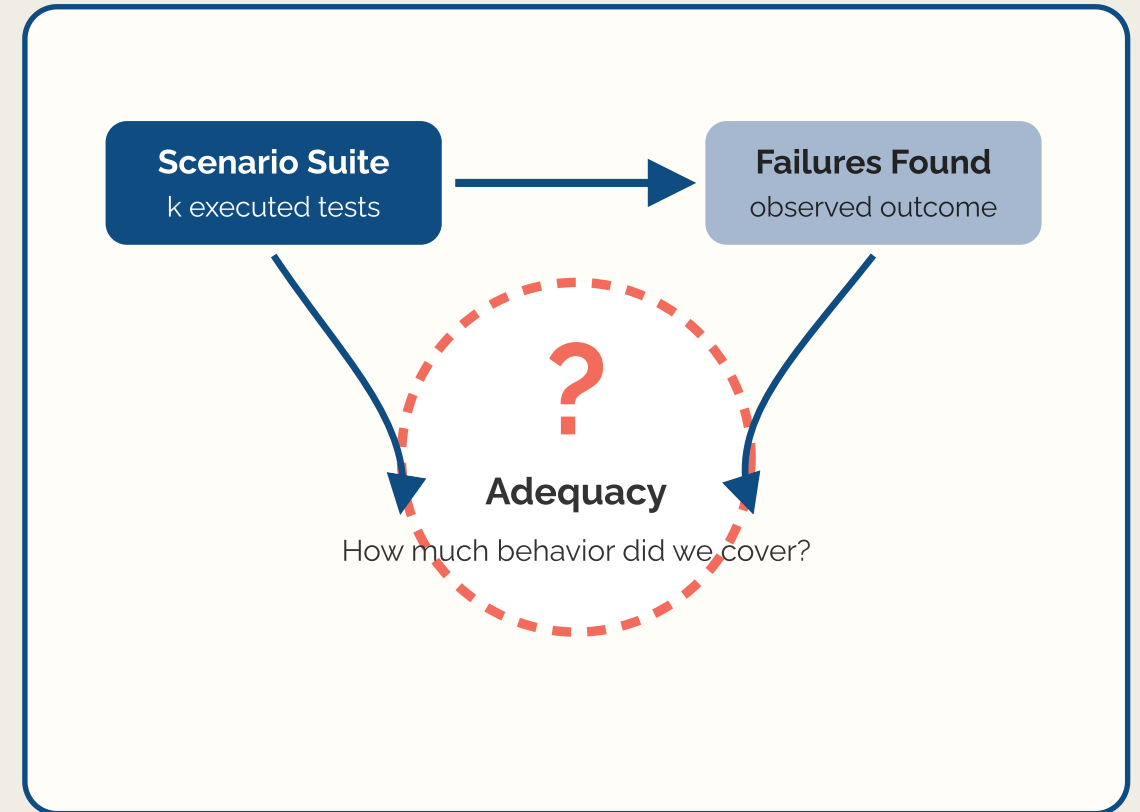
### Search

GA / RL / fuzzing

Scenario input diversity is structured, but still enormous.

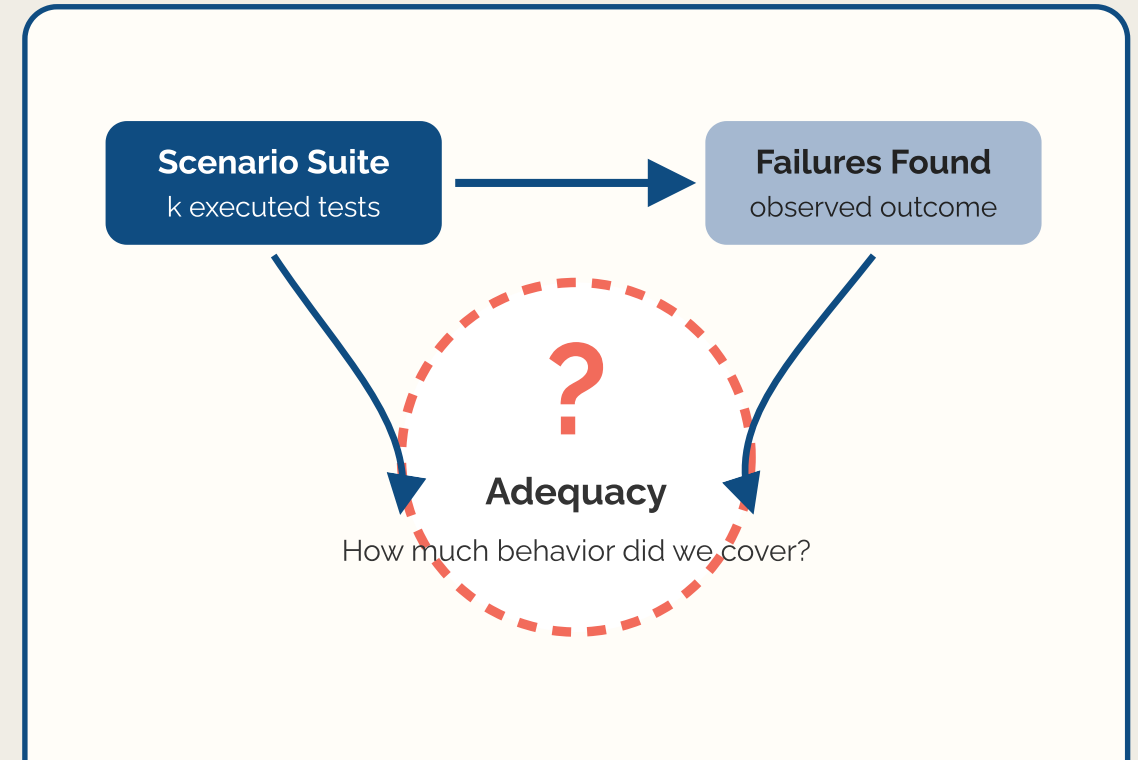
# Input Diversity Still Leaves A Gap

- **Critical-scenario search finds bugs efficiently**
- **But it rarely answers: when is a suite adequate?**
- **Scenario coverage can ignore ADS-internal behavior**
- **More scenarios can still be behaviorally redundant**



# Input Diversity Still Leaves A Gap

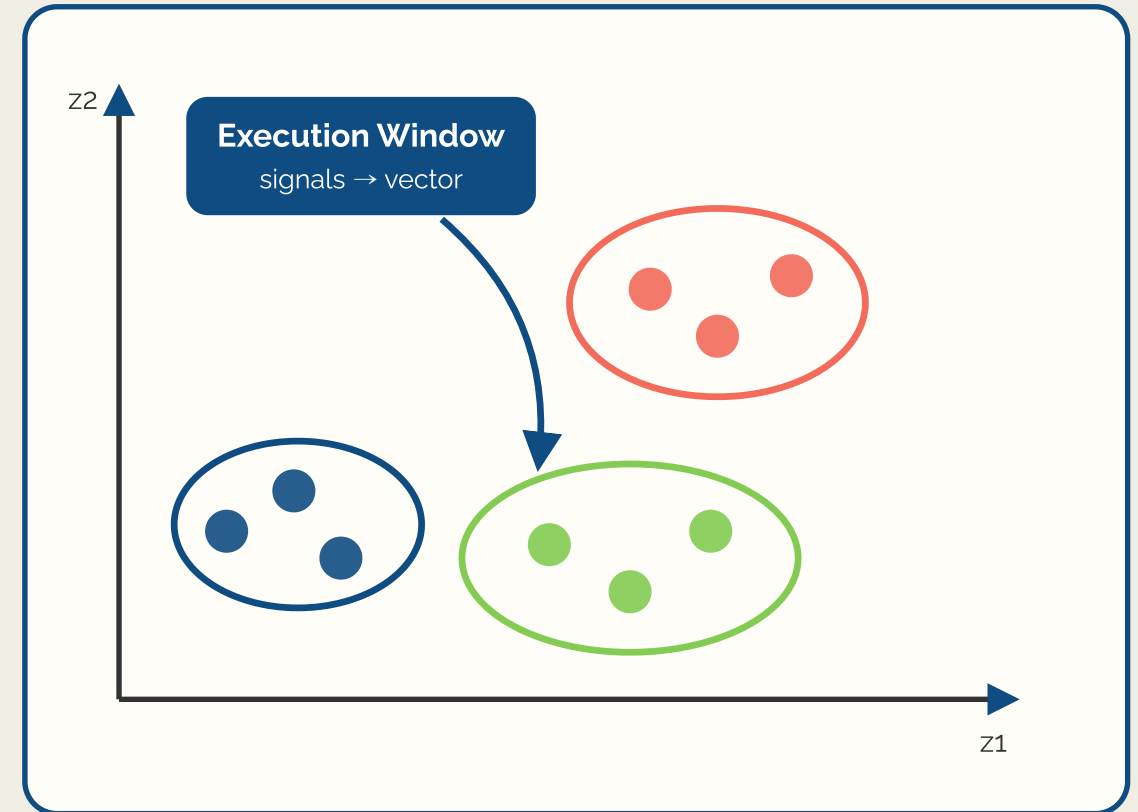
- Critical-scenario search finds bugs efficiently
- But it rarely answers: when is a suite adequate?
- Scenario coverage can ignore ADS-internal behavior
- More scenarios can still be behaviorally redundant



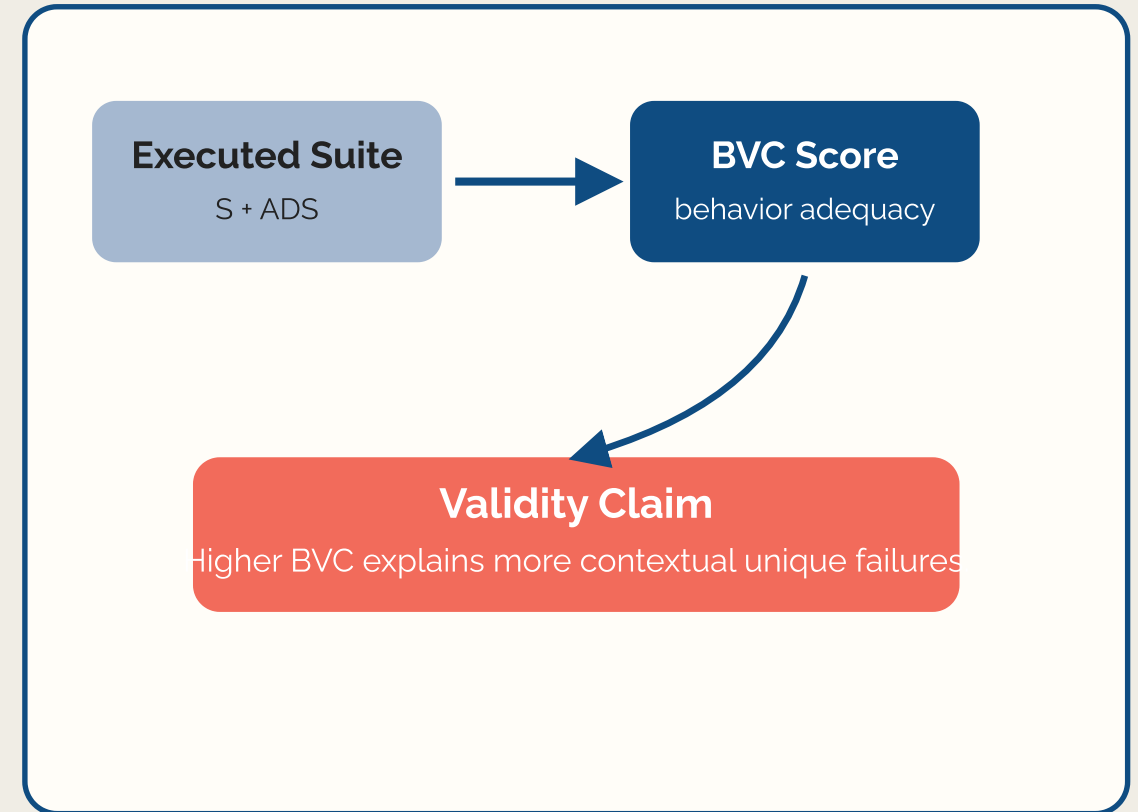
**We need a stopping/assessment signal after execution.**

# Behavior Vector Space

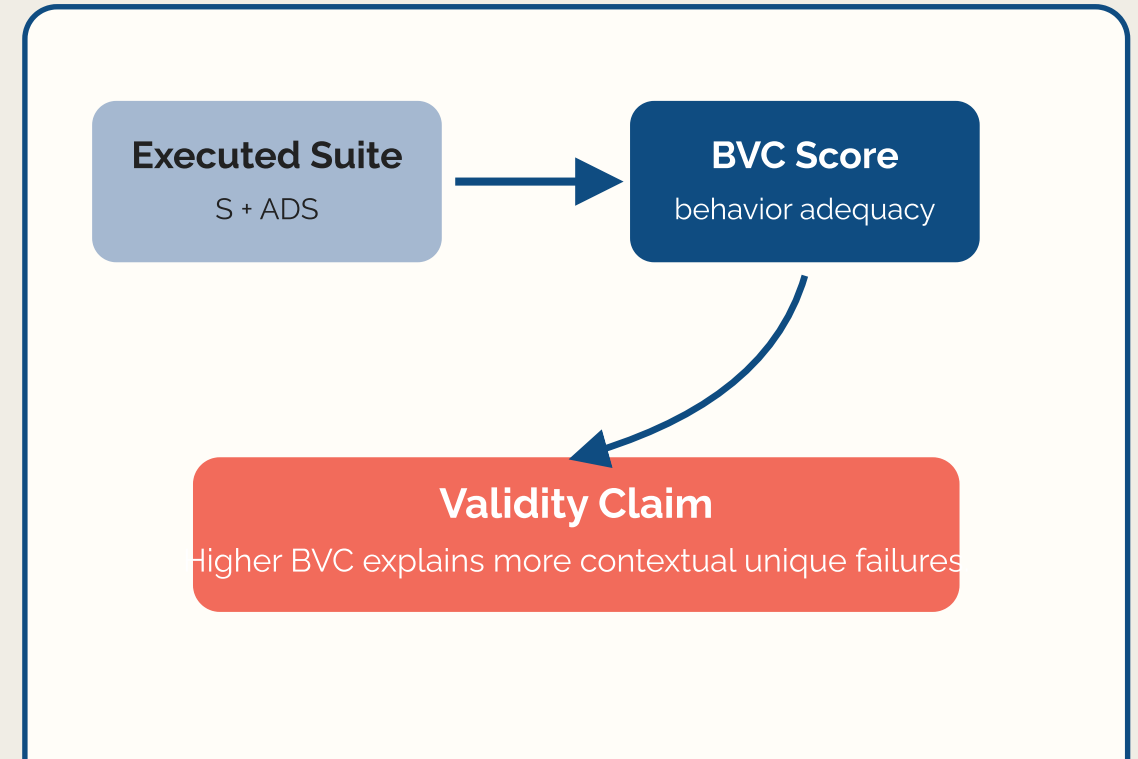
- Collect signals during each execution window
- Map each window to a normalized behavior vector
- Define coverage over the observed vector universe
- Keep the claim empirical, not global reconstruction



- **Input: executed scenario suite and ADS implementation**
- **Output: behavior-vector coverage and validity evidence**
- **Question: did the suite exercise diverse ADS behavior?**
- **Evidence: higher BVC should explain unique failures**



- **Input: executed scenario suite and ADS implementation**
- **Output: behavior-vector coverage and validity evidence**
- **Question: did the suite exercise diverse ADS behavior?**
- **Evidence: higher BVC should explain unique failures**

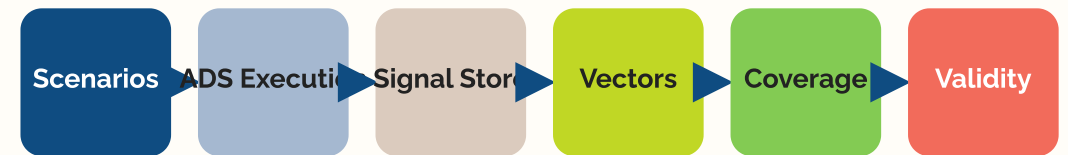


**Core thesis: BVC estimates how much observed ADS behavior a suite covers.**

# Overall Approach

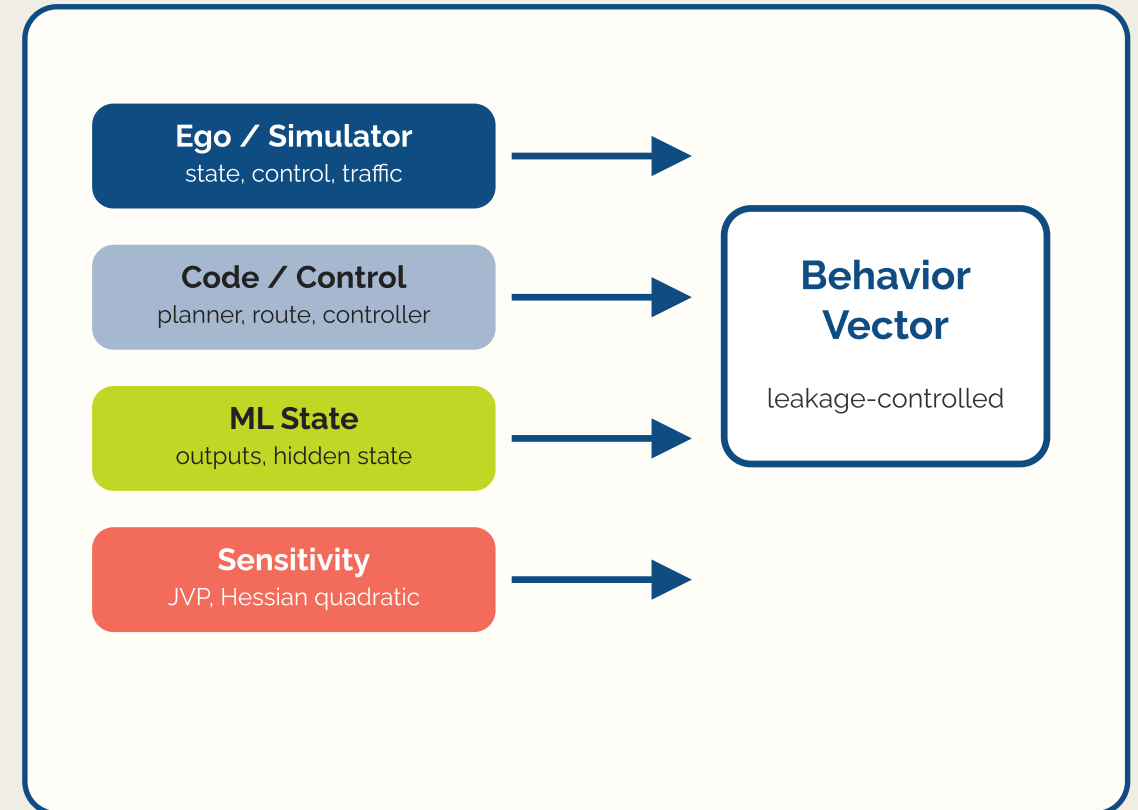
- Execute scenarios on a fixed ADS
- Record simulator, code, ML, sensitivity signals
- Aggregate windows into behavior vectors
- Score suite coverage over observed behavior
- Validate against unique failure discovery

## Post-Test Adequacy Pipeline



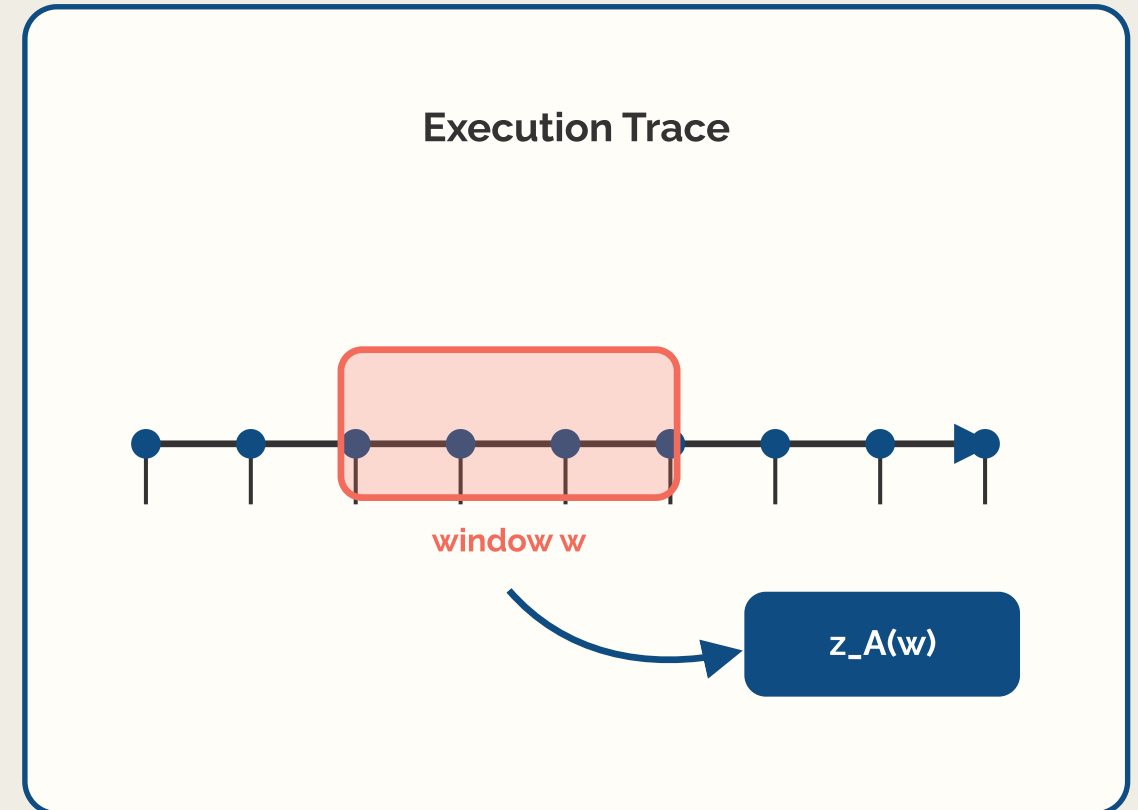
Coverage is computed after the suite has been executed.

- **Ego/simulator: state, control, traffic context**
- **Code/control: planner and controller decisions**
- **ML state: outputs, hidden representations**
- **Sensitivity: JVP and Hessian-quadratic probes**
- **Leakage control removes shortcut failure proxies**



$$z_A(w) = T(\varphi_A(w)) \in \mathbb{R}^d$$

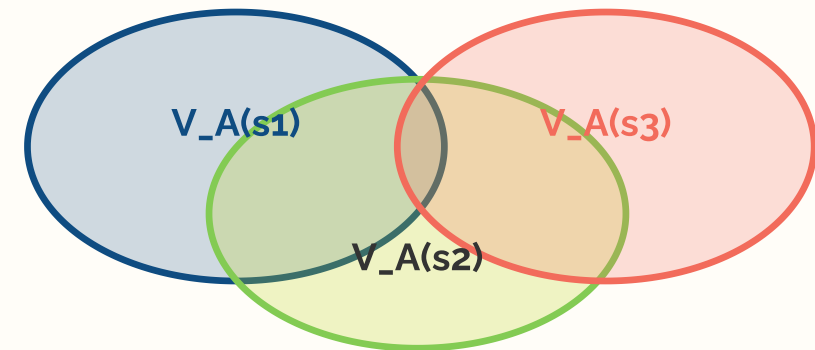
- **Windowed aggregation over execution traces**
- **Z-score normalization over the corpus**
- **Zero-variance removal and clipping**
- **Optional block-wise PCA for high-dimensional blocks**



$$V_A(S) = \bigcup_{s \in S} V_A(s)$$

$$U_A = \bigcup_{s \in D} V_A(s)$$

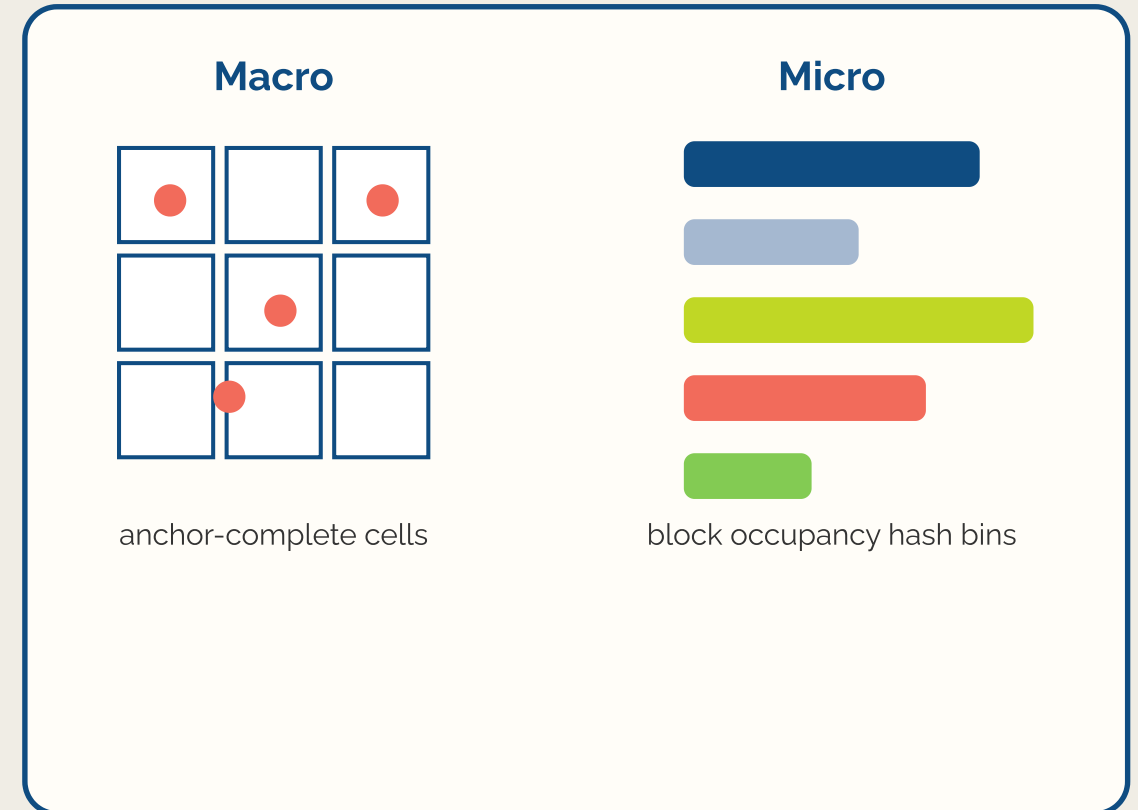
- Each scenario induces many behavior windows
- A suite is the union of its scenario vectors
- Coverage is measured against the executed pool



Suite and universe are finite unions of behavior vectors

# Macro and Micro Coverage

- **Macro: anchor-complete cells in full vector space**
- **Micro: block-level occupancy proxy**
- **Radius from nearest-other-scenario quantiles**
- **Main conservative metric: macro q25**
- **Robustness reports q05, q10, q25, q50**



RQ1. [Validity]

Is BVC positively associated with unique failure discovery?

RQ2. [Incremental Value]

Does BVC explain failures beyond scenario/input and simpler execution baselines?

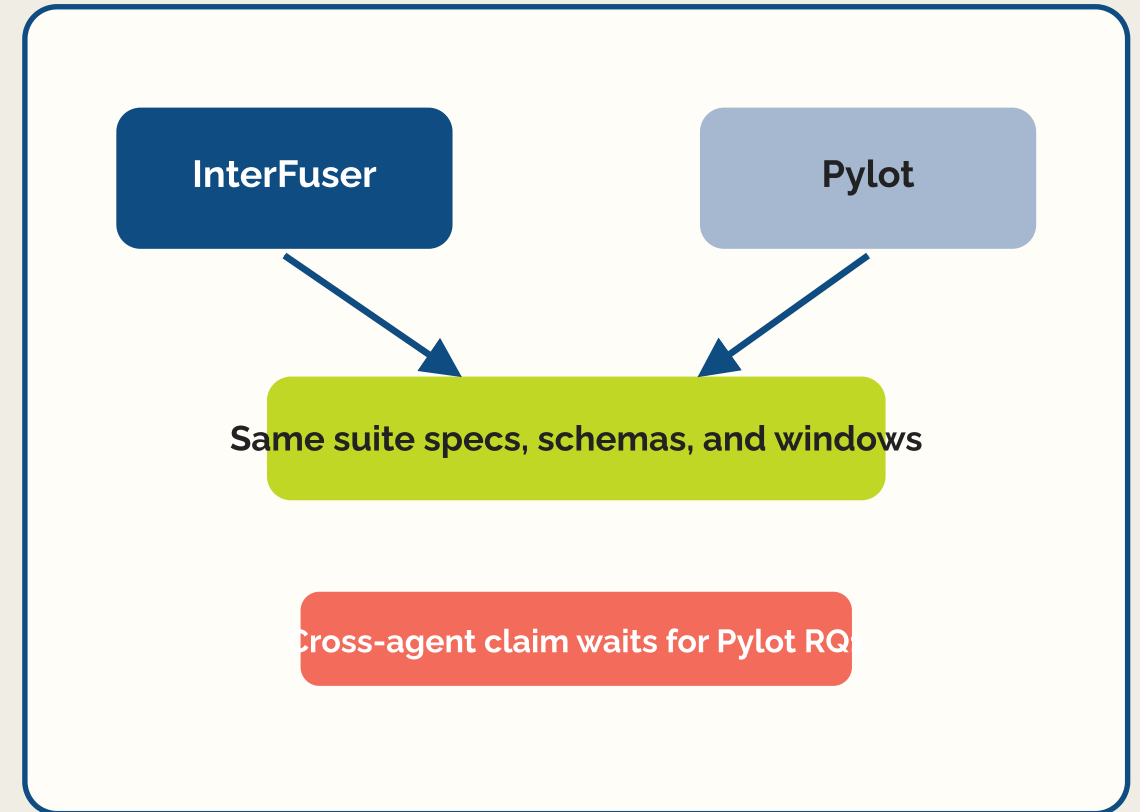
RQ3. [Construct Validity]

Which signal families make BVC useful, and do they align with failure semantics?

RQ4. [Robustness]

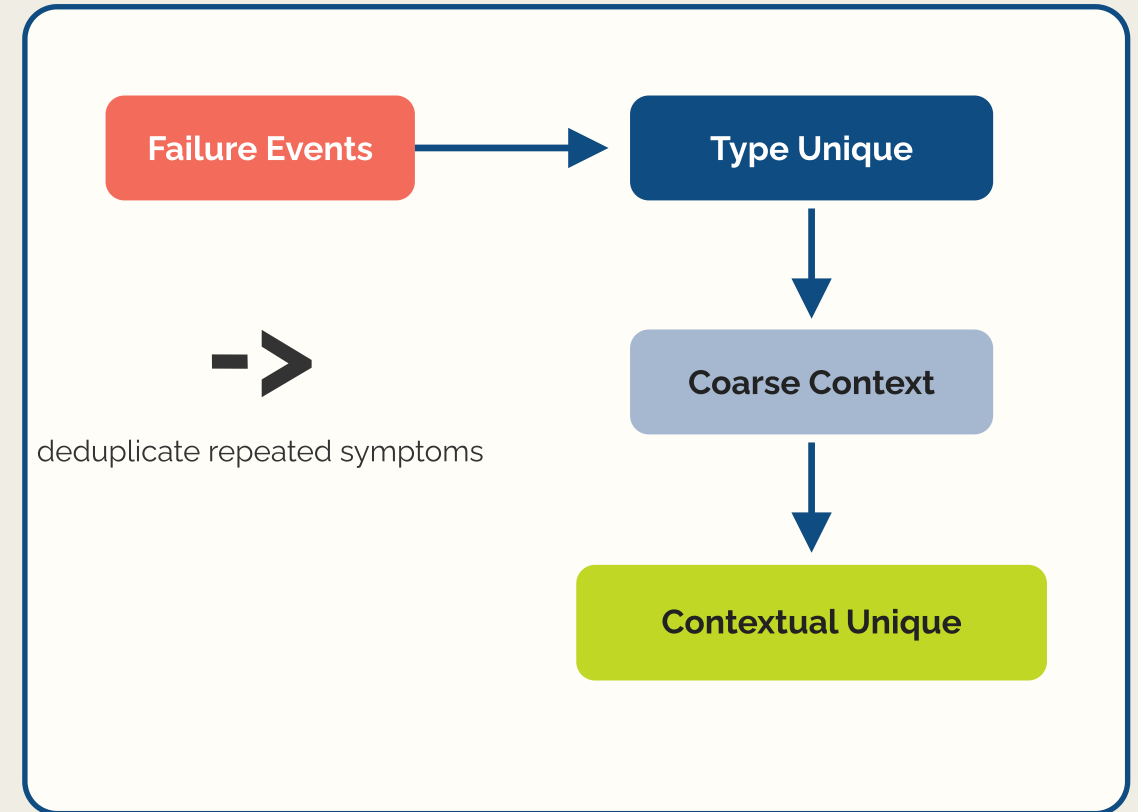
How sensitive and practical is BVC after signal collection?

Component	Current plan
ADS	InterFuser complete; Pylot running
Simulator	CARLA / Drivora scenario runner
Suite sizes	5, 10, 25, 35, 50, 75, 100
Suite specs	100 random suites per size
Windows	0.25s, 0.5s, 1.0s
Schemas	code_only, ml_only, all



# Unique Failure Definition

- **Failure event: one oracle/infraction occurrence**
- **Type-unique: failure family only**
- **Coarse unique: type plus coarse location**
- **Contextual unique: type plus scenario context**
- **Primary target: contextual unique failures**



- **Mode: adequacy\_fullset**
- **700 fixed suite specs reused everywhere**
- **3 schemas × 3 windows × 7 suite sizes**
- **Macro/micro metrics at q05, q10, q25, q50**
- **Correlation and OLS HC3 regression**

Primary choice	Value
Adequacy	macro_cell_coverage_q25
Failure target	contextual_unique_failure_count
Distance	L2 / sqrt(dim)
Controls	suite_vector_count
Rows	63 primary combinations

# RQ1 Dataset Snapshot

200

Successful scenarios

588

Ego traces

514

Failure events

249

Contextual unique

700

Fixed suites

1.38M

Vector rows

3,637

All normalized dim

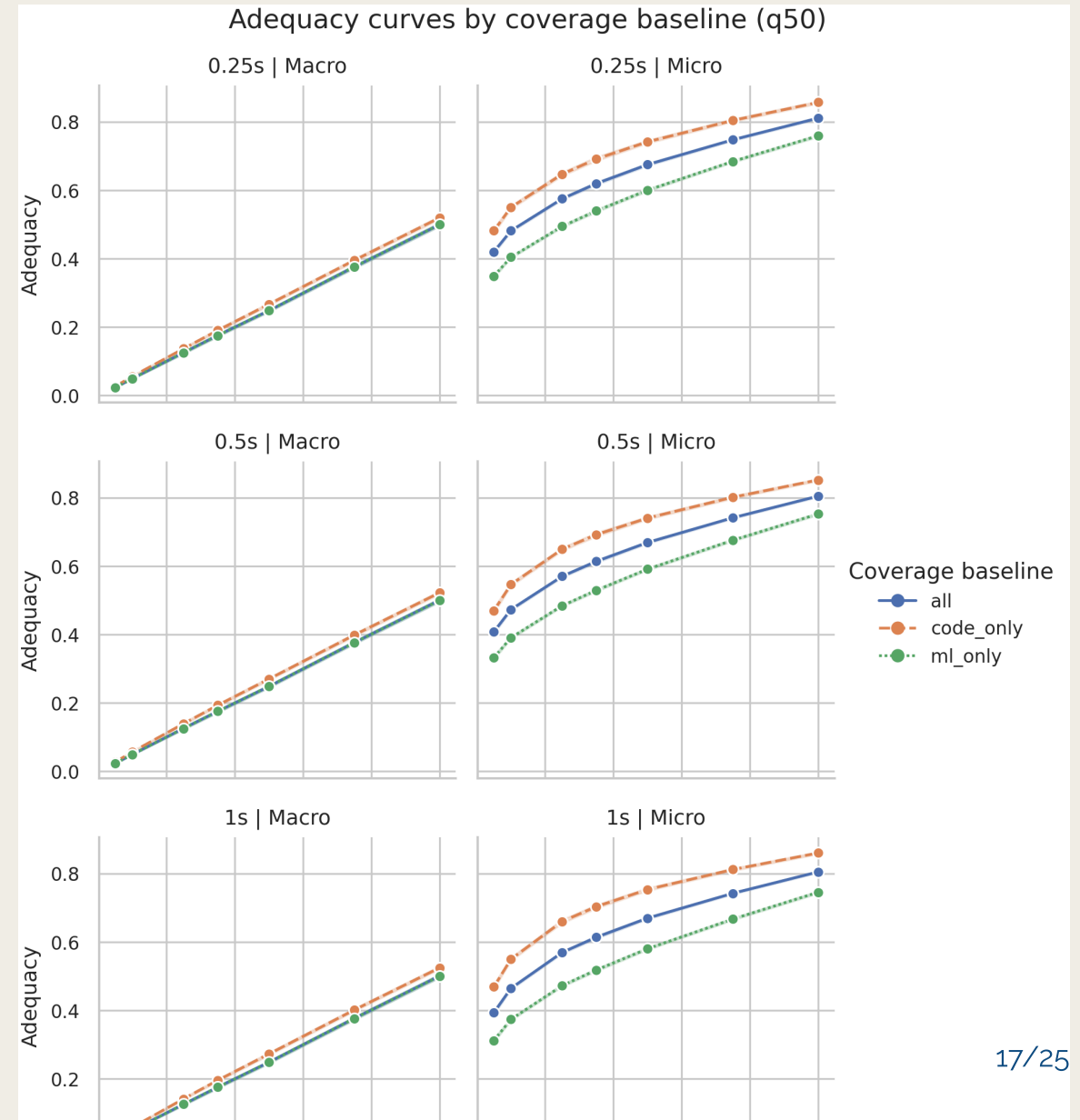
63

Primary rows

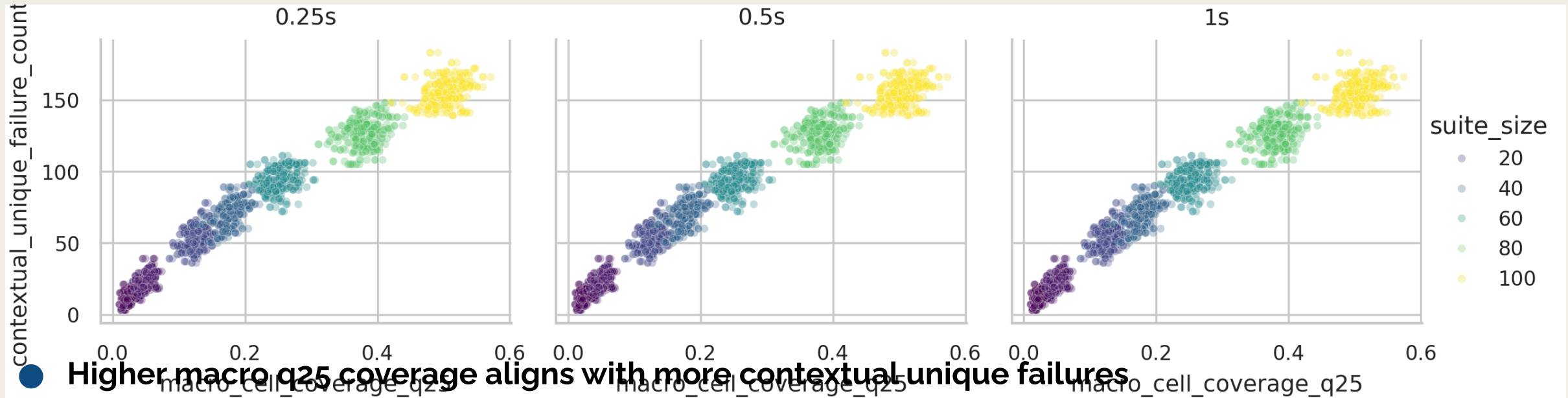
Source: [results/paper/tables/interfuser\\_rq1\\_dataset\\_summary.csv](#); [results/interfuser/rq/rq1/multi\\_window\\_summary.json](#)

# Adequacy Rises with Suite Size

- Coverage increases monotonically with suite size
- Macro grows conservatively from near zero
- Micro starts higher and saturates faster
- `code_only` is strong in macro q50
- `all` is strongest in micro q50



# Coverage Tracks Unique Failures



- Higher macro q25 coverage aligns with more contextual unique failures

- The pattern repeats across 0.25s, 0.5s, and 1.0s windows

- Suite size explains level, BVC explains ordering within size

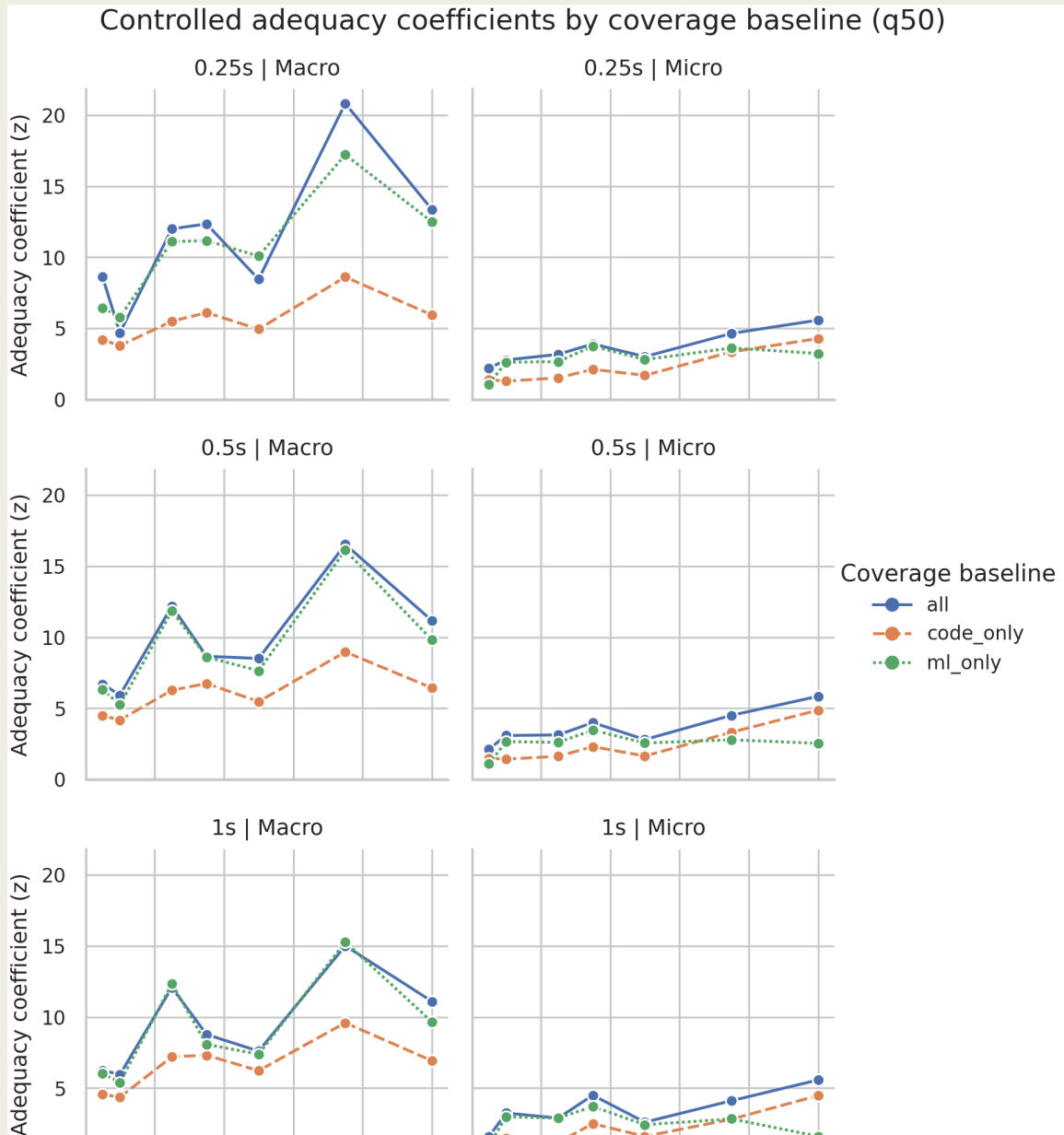
- This is the main validity picture

# Association Holds Across Sizes

- Primary Spearman: 63/63 positive
- Mean Spearman: 0.389
- Range: 0.232 to 0.555
- Positive in every schema-window-suite cell
- Strong support for InterFuser validity



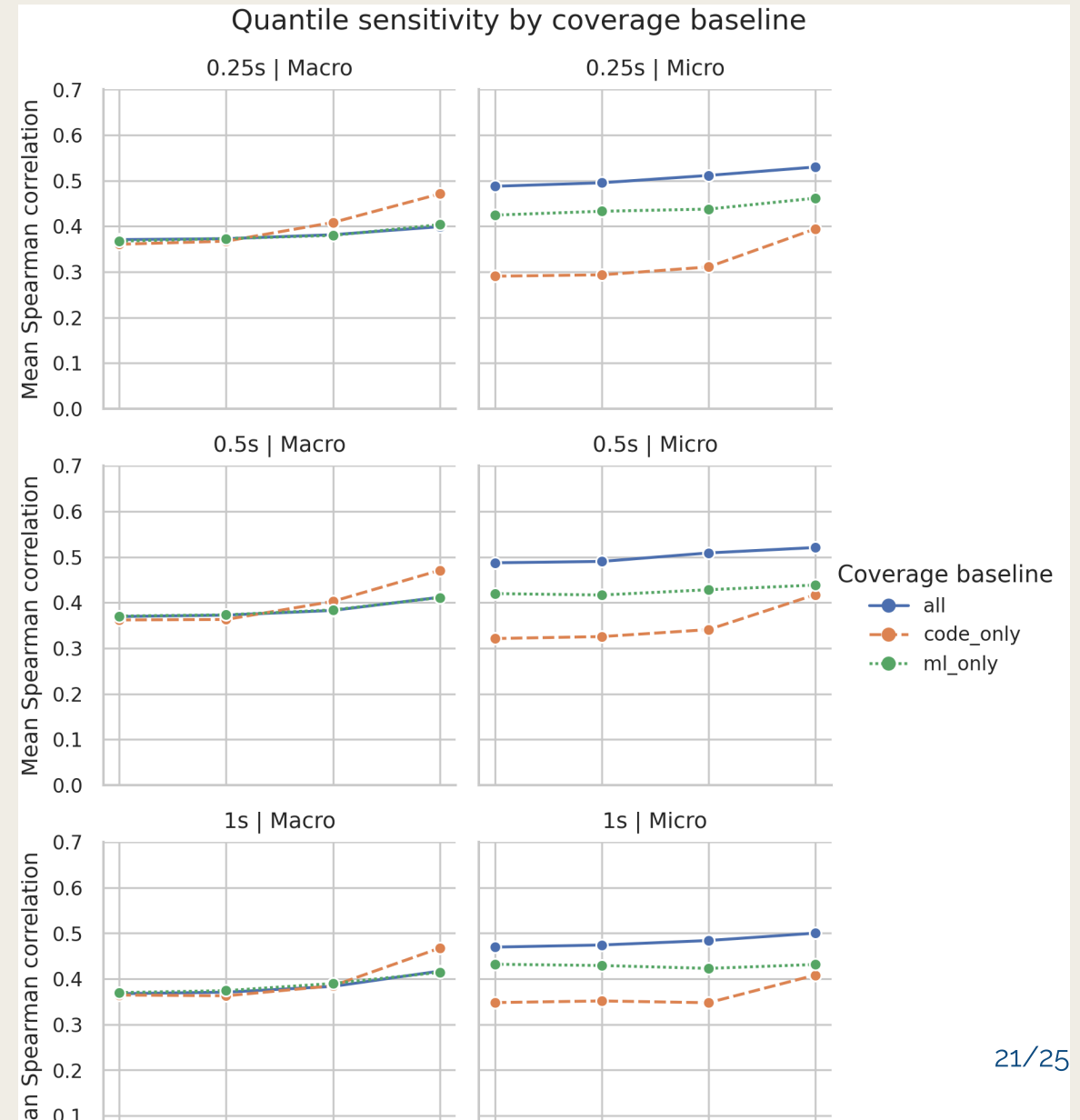
# Controlled Models Support BVC



- **Primary BVC coefficients: 63/63 positive**
- **Significant positive: 54/63 at  $p < 0.05$**
- **Significant negative: 0/63**
- **BVC remains explanatory after execution volume control**
- **Interpretation: not just more logged vectors**

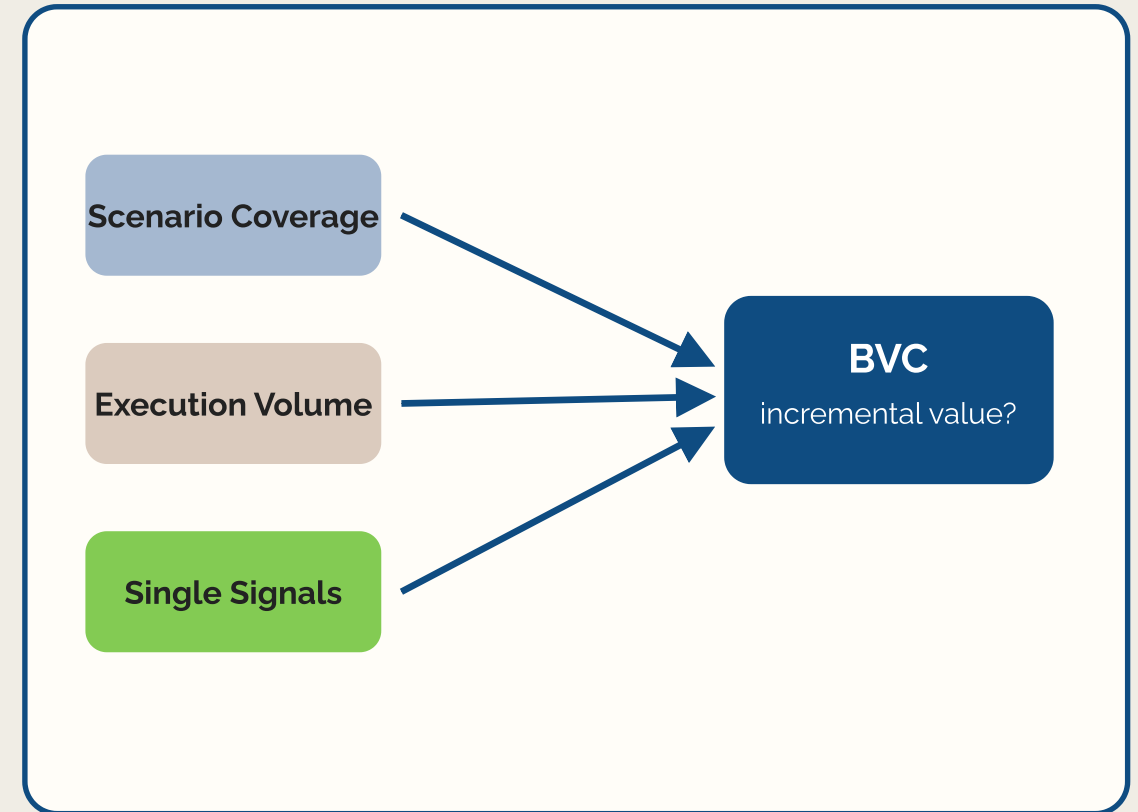
# Quantiles and Families Matter

- Macro q25 is the conservative primary metric
- q50 often gives stronger empirical association
- Micro/all is strong but less geometrically exact
- code\_only, ml\_only, all are separate baselines
- Do not collapse everything into one BVC score

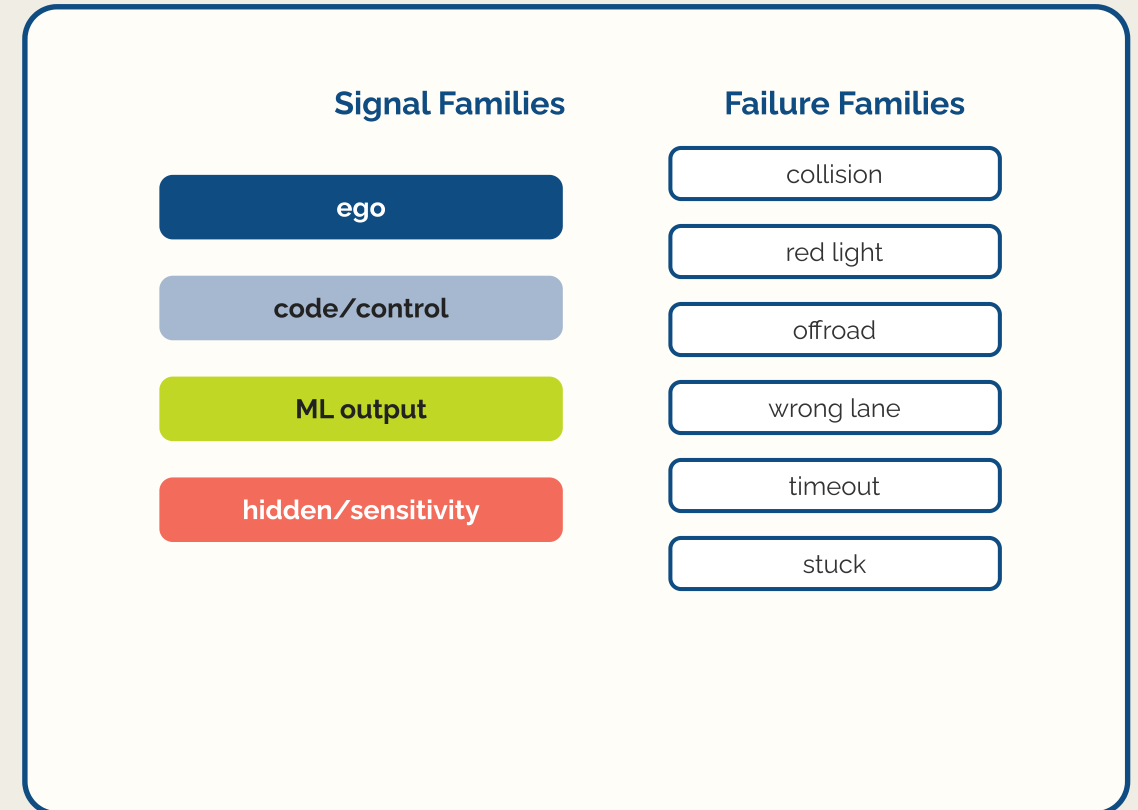


# RQ2 Baseline Comparison Plan

- Reuse RQ1 fixed suite specs
- Compare against random suite sampling
- Scenario/input-space coverage baselines
- Execution-volume and trace-count baselines
- Single-signal or family-only adequacy baselines



- Ablate signal families and schema variants
- Measure association with contextual failures
- Map signal families to failure families
- Use GA/search diagnostics as interim evidence
- Expected output: family × failure heatmap



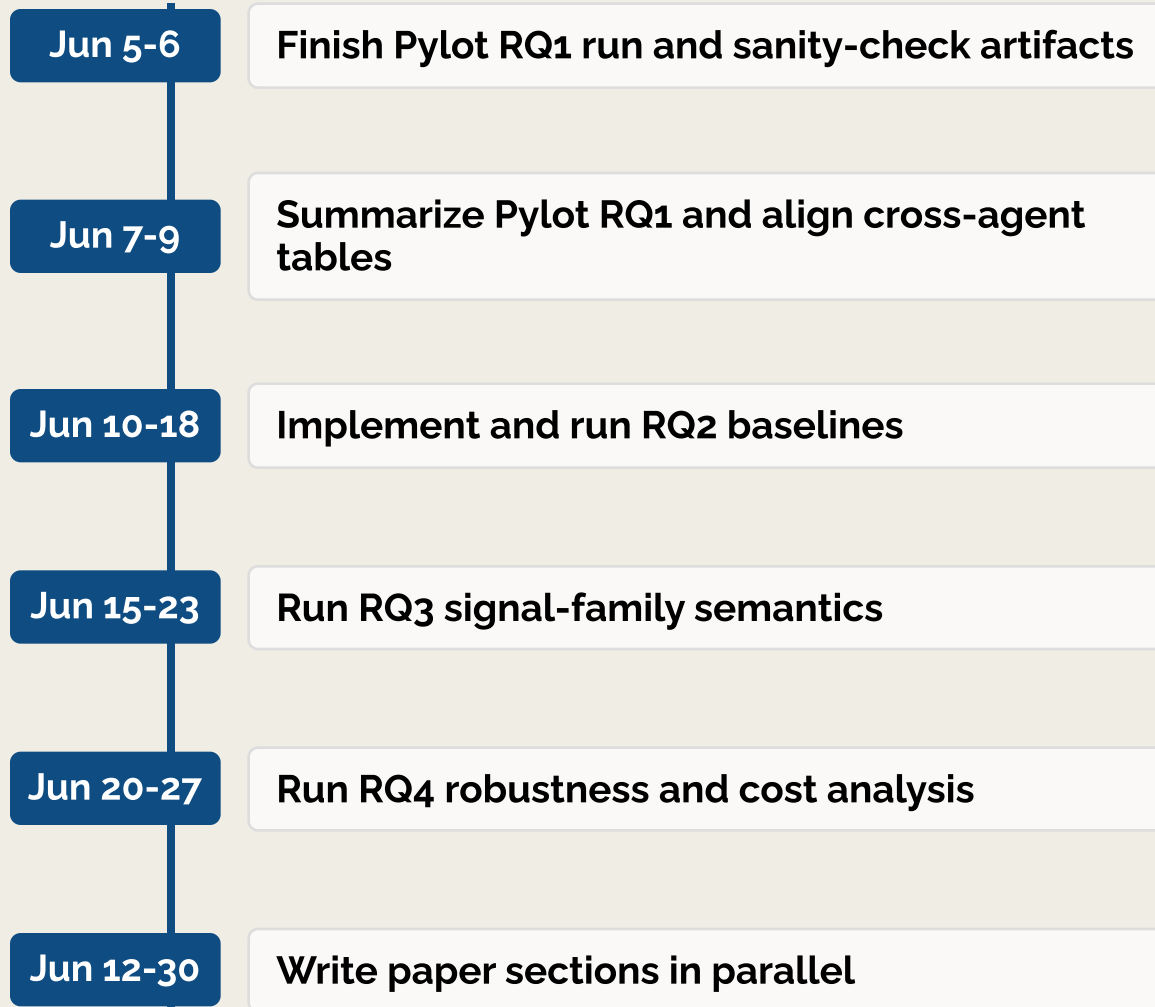
- Window sizes: 0.25s, 0.5s, 1.0s
- Radius quantiles: q05, q10, q25, q50
- Distance metrics: L2 now; cosine/group-weighted planned
- Failure granularities: event, type, coarse, contextual
- Cost: signal collection, vector build, cached scoring

## Robustness Grid

window	radius	distance	failure
0.25/0.5/1s	q05..q50	L2/cosine	context
cost	storage	GPU/CPU	cache

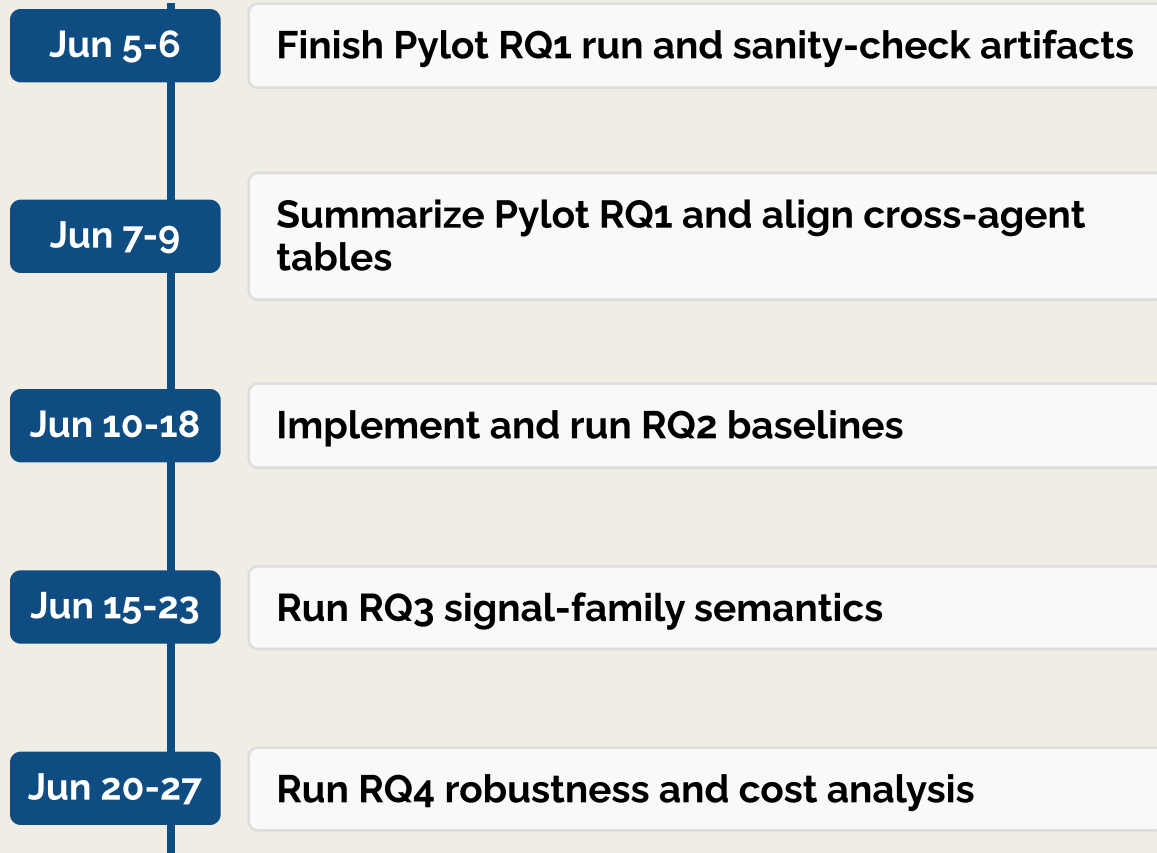
The claim should remain qualitatively stable across reasonable choices.

# TODO Through June 30



- **Pylot RQ1: in progress, target June 6**
- **RQ2/RQ3/RQ4 share RQ1 suite specs**
- **Summaries and figures must stay agent-scoped**
- **By June 30: full experiment narrative draft**

# TODO Through June 30



- **Pilot RQ1: in progress, target June 6**
- **RQ2/RQ3/RQ4 share RQ1 suite specs**
- **Summaries and figures must stay agent-scoped**
- **By June 30: full experiment narrative draft**

**Priority: finish cross-agent RQ1 first, then reuse caches aggressively.**

# Current Progress

Behavior Vector Coverage for ADS Test Adequacy

InterFuser RQ1 supports the core validity claim

{'Next checkpoint': 'Pylot RQ1 and baseline comparison'}

June 2026