

CERT: Certified Route Planning under Drifting Costs

Conformal certificates, sense-to-certify, and the price of staleness

Krishi Attri
Seoul National University
kattri@snu.ac.kr

June 10, 2026

Abstract

A scout robot routing through terrain whose costs drift — mud after rain, traffic after an incident — faces a question classical replanning never answers: *how good is the current route, given that most of the map is stale?* CERT answers it every round with a certificate: a high-probability bound $LB \leq OPT \leq UB$ on the optimal route cost, built from age-weighted non-exchangeable conformal prediction over drift-adjusted residuals, with paid sensing directed at the edges that shrink the certified gap fastest. We prove coverage at the claimed level with a staleness correction that degrades the claim visibly rather than silently; a certifiability *threshold* — a target gap is sustainable iff the sensing rate beats drift, so certification is a rate, not a state; a \sqrt{L} -tighter sum-aware upper certificate, including the selection-bias hazard it creates and the gate that controls it; and an impossibility theorem showing the certificate’s asymmetry is optimal. On replayed traffic from two cities the certificate holds even where real incidents violate the drift model up to half the time, and certificate-directed sensing achieves 2–3× lower travel-regret than freshness-, uncertainty-, or chance-driven sensing at equal budget.

1 Introduction

Figure 1 shows one round of the planner this paper builds and certifies.

Incremental replanners repair shortest paths quickly when the map changes [15, 19], informative path planners decide where to sense [25], and learned traversability models supply ever-better cost priors. What none of these provide is an online, sound answer to the operational question: *is the route I am about to execute within ϵ of the best currently-achievable route, and how confident am I, given that my last look at most of the map is minutes old?* The gap is not any single component. Measured suboptimality bounds exist for anytime search over *known* costs [18]; pay-to-sense edge resolution exists for *binary* blockages without certificates [3]; PAC path identification exists for *stationary* samplable costs [4]; and staleness-aware maps exist without route-quality guarantees [16, 26]. The intersection — a certificate on route cost that drives physical sensing, under continuously drifting, partially observed costs, maintained online — is, to our knowledge, unoccupied (Table 1 audits this claim family by family). The setting is the daily reality of deployed mobile robots: a scout vehicle whose off-road traversal costs change with weather and erosion, a delivery fleet whose street segments drift with traffic, an inspection UGV that must decide whether to commit to a corridor or spend another ranging scan first. In each case the operator’s question is not “what is the shortest path on my map?” but “can I trust the route my stale map implies, and if not, where do I look?” — which is precisely the certificate and the sensing rule CERT provides.

CERT occupies it with a deliberately asymmetric design. Each edge carries a point estimate, an observation age, and a drift-rate bound; conformal calibration over drift-adjusted residuals converts these into per-edge cost intervals whose coverage guarantee degrades *explicitly* with both the miscoverage target and a staleness correction. Two incremental searches maintain the optimistic and conservative shortest paths; their costs are the certificate. Sensing is pointed at the certified gap: the planner observes whichever route-critical edge buys the largest expected gap reduction per unit cost, with an age-triggered round-robin backstop that converts greed into a guarantee. The robot executes only conservative incumbents, and a conductivity

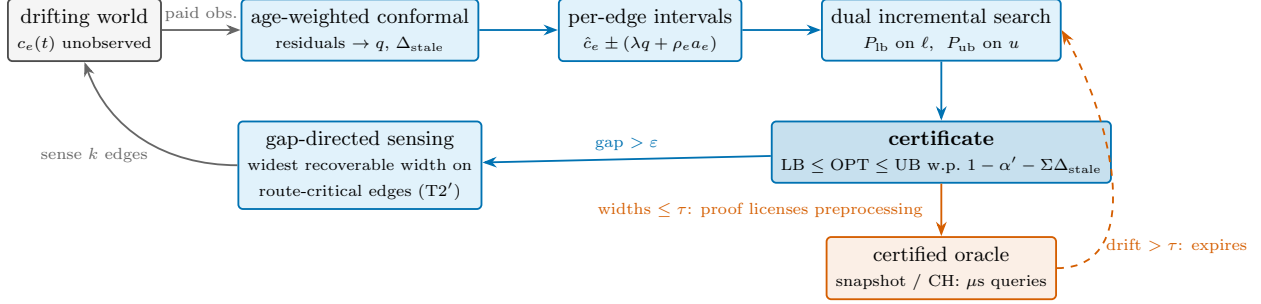


Figure 1: One CERT round. Paid observations feed an age-weighted non-exchangeable conformal scorer whose quantile prices noise (q) and staleness ($\rho_e a_e$) into per-edge intervals; two incremental searches bound the optimum from both sides, yielding a certificate that either certifies ($\text{gap} \leq \varepsilon$) or directs the next observations at the edges that shrink it fastest. When the certificate proves every interval within τ , that proof licenses lookup-speed preprocessed queries (Section 6.7) — revoked the moment drift exceeds τ .

hysteresis selects among certified-equivalent incumbents to suppress solution churn without touching the certificate.

The framework forced four findings we did not anticipate at design time. (i) *Certification is a rate*: a certificate self-extinguishes even in a static world unless maintenance sensing sustains it — the certifiability threshold (Theorem T2') makes the required rate explicit, including an impossibility direction. (ii) *Adaptive and static guarantees fight*: an adaptive conformal layer silently cancels a provable widened margin; the provable mode must freeze it. (iii) *Tighter bounds buy new hazards*: replacing the Bonferroni union bound with a sum-aware block quantile (Theorem T4) tightens the gap by $\Theta(\sqrt{L})$ but exposes the optimizer's winner's curse — measured as a coverage drop from 1.000 to 0.823 — which a freshness-gate protocol controls. (iv) *Locality and staleness conflict*: age-driven widening touches every edge every round and destroys incremental-search locality; a lazy pre-widening cache restores it with a soundness proof and a quantified width price.

Contributions. (1) The problem formulation and certificate construction (Sections 2–3). (2) Theorems T1a/T1b (coverage with explicit staleness correction; observable vs. latent semantics), T2' (certifiability threshold, both directions), T4 (sum-aware upper certificate) with the freshness-gate protocol, and the asymmetry theorem T5 (Section 4; full statements and proofs in Appendix A). (3) Certificate-gated preprocessing: a *proof*, rather than an assumption, licenses lookup-speed query structures, which expire automatically when drift exceeds their validity window — this is what lets a certified planner reach the microsecond query class of static-known route planners (Section 6). (4) The maintenance-sensing and lazy pre-widening mechanisms. (5) A four-tier empirical study with coverage validated against ground truth in simulation and on replayed real traffic, including misspecification, off-model, and boundary stress (Section 6).

2 Problem Setting and Assumptions

A directed graph $G = (V, E)$, start s , goal g ; unknown time-varying costs $c_e(t) > 0$; one paid noisy observation per round of length Δ . The planner must report, every round, a certificate (LB, UB, confidence) with $\text{LB} \leq \text{OPT}(t) \leq \text{UB}$ at the stated confidence, execute only conservatively validated routes, and stop sensing when $\text{UB} - \text{LB} \leq \varepsilon$.

Assumptions (stated fully in Appendix A; all are swept or stress-tested in Section 6): **A1** bounded drift $|c_e(t') - c_e(t)| \leq \rho_e |t' - t|$ with ρ_e conservatively known; **A2** the observation-residual distribution drifts in total variation at rate at most ε_{TV} ; **A3** symmetric unimodal observation noise (needed only for the latent-cost and sum-aware variants); **A4** edges sharing a calibration buffer share a noise family.

3 The CERT Planner

Certificate substrate. Re-observing edge e at age a yields the drift-adjusted score $R = |Y - \hat{c}_e| - \rho_e a$; scores enter a rolling buffer with age-geometric weights $w_i = \rho_w^{\text{age}_i}$ (data-independent, as the non-exchangeable guarantee requires [1]). The weighted $(1-\alpha)$ quantile q gives intervals $\hat{c}_e \pm (\lambda q + \rho_e a_e(t))$; $\lambda = 1$ certifies observables (T1a), $\lambda = 2$ certifies latent costs (T1b). The reported confidence is $1 - \alpha' - \sum_e \Delta_{\text{stale}}$ — staleness degrades the claim visibly. An adaptive-conformal tracker [12] supplies an assumption-free long-run safety net; it is frozen in the provable mode (Section 6.1, ACI-interaction finding).

Dual incremental search. Two D* Lite instances maintain the ℓ -shortest and u -shortest paths; LB is the ℓ -cost of the former, UB the best u -cost among candidates (any path’s u -cost upper-bounds OPT on the coverage event). Age widening would touch every edge every round; a *lazy pre-widening* cache computes metrics at age $+B\Delta$ so entries stay conservatively valid for B rounds, restoring repair locality at width price $2\rho_e B\Delta$ per edge (Lemma A.5, Appendix A.5).

Sense-to-certify. Since $\text{UB} - \text{LB} \leq \sum_{e \in P_{\text{lb}}} (u_e - \ell_e)$, sensing only the optimistic path’s edges suffices to control the gap. The selector takes the route-critical edge maximizing expected gap recovery ($2\rho_e a_e$) per sensing cost; an age-triggered backstop forces round-robin re-observation, which is what makes the achievability theorem apply to the deployed policy. Certified rounds still sense on projected certificate expiry and a calibration-freshness floor (maintenance sensing); without it the claim self-extinguishes even in static worlds.

Hysteresis and the freshness gate. Among incumbents whose u -cost is within a slack of UB, a decaying edge-conductivity κ selects the stickiest — churn drops 70% at zero certificate cost (the bound reported is unchanged; only the executed selection moves). The same stability opens the *freshness gate* for the sum-aware upper certificate (T4): the \sqrt{L} -margin bound applies only when every incumbent edge has been re-observed since the path became incumbent, which restores the fixed-path premise that the optimizer’s selection otherwise violates.

4 Theory

The full statements and proofs are given in Appendix A; we summarize the main results here. **T1a/T1b (coverage):** per-edge intervals cover the next observation at level $1 - \alpha - \Delta_{\text{stale}}$ under A1–A2, and the latent cost at a doubled margin and level $1 - 2\alpha - 2\Delta_{\text{stale}} - \pi_{\text{cal}}$ under A1–A3 via Anderson domination; Δ_{stale} is computed from realized calibration ages, sharpening the age-uniform corollary of Barber et al. [1]. **T2' (certifiability threshold):** round-robin sensing of the L optimistic-path edges sustains $\text{UB} - \text{LB} \leq 2L\bar{q} + \bar{\rho}\Delta L(L - 1)$; conversely no policy sustains $\varepsilon < 2mq_{\min} + \rho_{\min}\Delta m(m - 1)$ across an m -edge cut, and $\varepsilon < 2Lq$ is unattainable at any sensing rate. Static worlds and the deterministic scout stopping rule of Rockenbauer et al. [25] are the $\rho \rightarrow 0$ and $q \rightarrow 0$ corners. **T4 (sum-aware upper certificate):** block-conformal calibration of signed deviation sums gives a fixed-path upper bound with $\Theta(\sqrt{L})$ margin replacing Bonferroni’s $\Theta(Lq_{\alpha'/L})$; the certificate becomes asymmetric because the lower bound must hold uniformly over paths. The fixed-path premise fails for optimizer-selected incumbents (winner’s curse) and is restored conditionally by the freshness gate. **T5 (LB impossibility):** the asymmetry is forced — on layered graphs any valid uniform lower bound pays $\Omega(L\sigma\sqrt{\ln w})$ slack (posterior-greedy selection over exponentially many paths), which per-edge Bonferroni matches up to logarithmic factors. **T6 (decision-uniform validity):** per-round claims are marginal; a policy that acts on certificates selects rounds, and α -spending over the (few) decision instants restores trajectory-level validity exactly where the trajectory consumes it — while full per-round time-uniformity is quantifiably impractical ($n \gtrsim 63\text{k}$ calibration scores at Bonferroni levels). **T7 (churn-measured floor):** under drift the optimistic path hops over a churn set of $K \geq L$ edges; the certifiability floor and the adaptive sensing rate must use the online-tracked \hat{K} , and focused sensing *suppresses* churn ($K: 59 \rightarrow 11 \approx L$ measured) rather than chasing it.

5 Related work

CERT sits at the intersection of four research lines that have not previously been combined: search with cost certificates, sensing allocated to resolve path decisions, conformal prediction for graph costs, and temporal map models. We organize the discussion by these lines and state, for each, the precise position CERT occupies. We do not claim to be the first certified planner, the first active-sensing replanner, the first staleness model, or the first bio-inspired memory; the contribution is the conjunction, together with the staleness correction Δ_{stale} and the certifiability threshold $T2'$ as new analytical objects.

Certified and bounded search. TASP [31] computes explicit lower and upper bounds on path cost from bounded estimators, but its bounds tighten through additional computation over static costs, with no sensing. Anytime search such as ARA* [18] and AD* [19] reports a measured suboptimality certificate, yet assumes known costs and closes the optimality gap by spending more search effort rather than by observing the world. Robust interval shortest-path formulations carry cost intervals through the search but treat the intervals as given and fixed. CERT differs on the source and target of the bound: its certificate is a high-probability statement about an *unknown, drifting* cost, and the gap is closed by paid observation rather than by computation.

Sensing to resolve a path. The Canadian Traveller Problem with remote sensing [3] pays to sense edges and places sensing by value of information, but optimizes expected cost over binary blockages and produces no cost certificate and no drift model. Edge-evaluation planners — BISECT, LazySP, and Generalized Lazy Search [5, 6, 21] — decide which edge to evaluate near-optimally, but evaluation is a computational query returning binary validity, not a cost interval. PAC combinatorial pure exploration [4] returns an (ε, δ) -certified best path by sampling arms, and InfoBAX [23] senses to identify a function property such as a shortest path; both assume a stationary problem, neither trades sensing cost against travel, and neither maintains an execution loop. CERT inherits the "which edge to sense" question from this line but answers it against a different objective: shrink a probabilistic cost *certificate* on a drifting graph while executing.

Scout and informative path planning with guarantees. Traversing Mars [25] scouts until a path is proven optimal or infeasible, using a deterministic stopping rule with no probabilistic bound, no staleness, and no incremental reuse; CERT recovers that rule as the degenerate case $\rho \rightarrow 0, q \rightarrow 0$ of $T2'$, and so strictly generalizes it. Informative path planning with guaranteed estimation uncertainty [14] certifies the variance of a sensed field uniformly over space; CERT instead certifies the route decision, which depends on cost differences along candidate paths rather than on the field everywhere.

Temporal and staleness map models. FreMEn [16] models periodic environment dynamics in the frequency domain, the persistence filter [26] estimates how long a feature remains valid, and BRULE [10] and persistent monitoring [27] treat revisitation as a resource over time-varying state. These works represent staleness as map or feature uncertainty. None couples observation age to an inflation of a cost interval and propagates that inflation into a route certificate, which is the mechanism Δ_{stale} and the $\rho_e a_e$ widening provide.

Conformal prediction in robotics. Conformal prediction beyond exchangeability [1] supplies the non-exchangeable coverage bound that Δ_{stale} instantiates, and adaptive conformal inference [12] supplies the assumption-free long-run safety net. Perceive with Confidence [22] applies conformal prediction to perception for safe planning, but on detections rather than edge costs. Two adjacent works are the closest neighbors. Luo and Zhou [20] build conformal intervals directly on sums of edge labels — path cost — through a sum-aware nonconformity score, avoiding the Bonferroni penalty, but only under exchangeability, with no drift, weights, sensing, or online loop. CQR-GAE [29] produces conformal edge intervals for robust shortest path, but assumes exchangeability, runs one-shot, and propagates no path-level coverage and no sensing. CERT uses the non-exchangeable bound to handle drift, optionally adopts the sum-aware score under those weights (the $T4$ upper certificate), and closes the loop with sensing and incremental maintenance.

prior family	drift-aware intervals	path-cost certificate	online incremental	gap-directed paid sensing
D* Lite / AD* [15, 19]	×	×	✓	×
LazySP / GLS [6, 21]	×	×	✓	×
CTP + sensing [3]	×	×	×	✓
Conformal sums (CIA) [20]	×	✓	×	×
CQR-GAE [29]	×	✓	×	×
TASP [31]	✓	×	×	×
FreMEn [16]	✓	×	×	×
IPP / BAX [14, 23, 25]	×	×	×	✓
E-Graphs / stigmergy [24, 9]	×	×	✓	×
CERT (this work)	✓	✓	✓	✓

Table 1: Each prior family misses at least two of the four properties whose conjunction CERT occupies. “Drift-aware intervals”: uncertainty that grows with observation age under an explicit drift model (exchangeable conformal constructions do not qualify — Figure 5 measures the consequence). “Path-cost certificate”: a coverage guarantee on $LB \leq OPT \leq UB$, not a heuristic bound on point estimates.

Memory and experience reuse. Experience graphs [24] reuse prior solution paths to speed search, and ACO and Physarum models [9, 30] reinforce decaying edge scalars as a search heuristic. CERT’s optional conductivity module κ occupies a narrow delta: it applies observation-age decay to warm-start a guaranteed incremental search, and is held outside the certificate by construction, so it can never affect coverage. This module is demoted to an ablation.

The unoccupied intersection. Table 1 makes the claim auditable: each neighbor misses at least two of the four properties {drift-aware weighted intervals, path-level coverage propagation, online incremental maintenance, sensing allocated to the certified gap}. The cell CERT occupies is their conjunction: non-exchangeable, age-weighted conformal edge-cost intervals under an explicit drift model, propagated to a path-cost certificate $LB \leq OPT \leq UB$, maintained online by incremental D* Lite, with sensing pointed at shrinking the certified gap and a threshold that states when no target gap is sustainable.

6 Experiments

The experiments follow the claim ladder. We first establish that the certificate covers the true optimum at the claimed rate and that the certifiability threshold is visible (Section 6.1). We then tighten the bound with the sum-aware upper certificate (Section 6.2) and audit its calibrated building block in isolation (Section 6.3). With the certificate validated, we measure travel regret against sensing baselines in unknown terrain (Section 6.4), confirm coverage on replayed real traffic (Section 6.5), and locate the claim’s boundaries on standard pathfinding maps and at road scale, where certificate-gated preprocessing reaches the microsecond query class (Sections 6.6–6.7). Incremental-repair latency, the component ablations, lifelong operation, and the churn-measured floor close the section (Sections 6.8–6.11). Coverage is a model-conditional claim, validated only in simulation and on recordings where an oracle runs Dijkstra on the true costs every round; reported coverage is the empirical rate among valid rounds, with Clopper–Pearson intervals. Throughout, results are stated in the present tense and numbers are medians over the seeds and rounds named in each caption. Every experiment below is regenerated by a single driver script under `scripts/` (e.g. `run_tier0.py`, `run_tier2.py`, `run_metr_la.py`, `run_movingai.py`, `run_ablations.py`, `run_lifelong.py`, `run_gaussian_break.py`), with per-experiment findings logged under `docs/results/`.

6.1 Tier 0: coverage and the certifiability threshold

Table 2 reports 25 seeds \times 300 rounds per condition on 6×6 grids with $\varepsilon = 5$, $\alpha' = 0.2$, $\varepsilon_{TV} = 10^{-4}$ unless noted, with maintenance sensing and lazy pre-widening enabled. Coverage is measured against the claimed line $1 - \alpha' - \sum \Delta_{\text{stale}}$. Figure 2 shows one representative gap trajectory and Figure 3 the coverage-versus-claim summary.

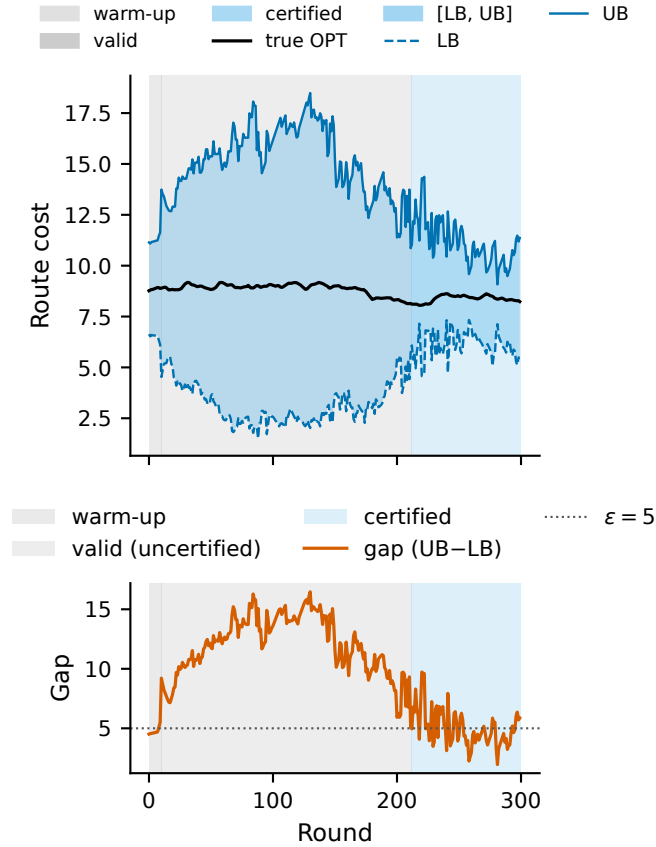


Figure 2: Certificate gap trajectory for one representative episode (6×6 grid, bounded drift $\rho=0.02$, `seed=3`, default CERT planner). *Top:* true optimal cost OPT (black) overlaid on the [LB, UB] certificate band (blue fill). *Bottom:* certificate gap UB-LB (red) and the target $\epsilon=5$ (dotted). Shading: grey = warm-up (annealing claim below target), mid-grey = valid-uncertified, blue = certified ($\text{gap} \leq \epsilon$). Single episode shown; aggregate over 25 seeds \times 300 rounds is reported in Table 2.

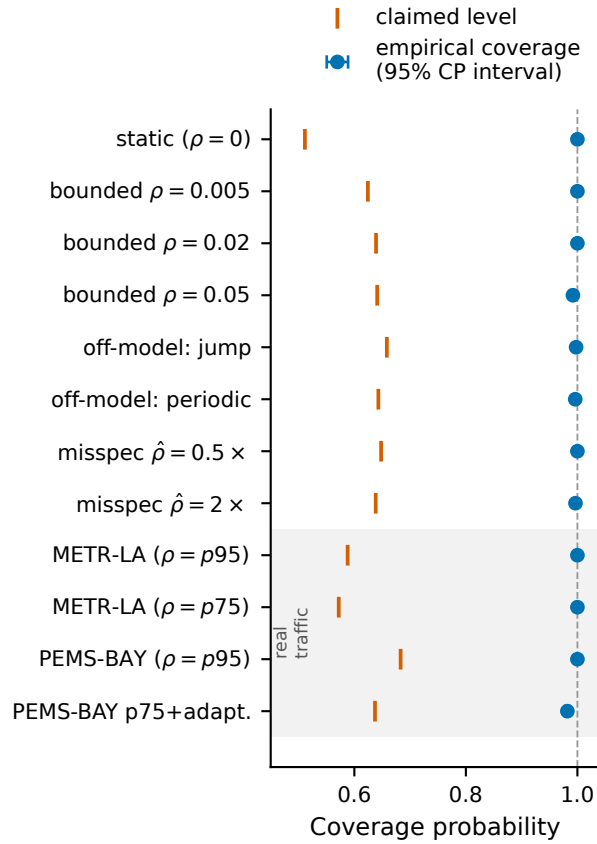


Figure 3: Empirical coverage (filled circles, 95% Clopper–Pearson bars) vs claimed confidence level (red tick marks) for eight synthetic conditions (white background) and four real-traffic conditions (grey background; METR-LA 20 days, PEMS-BAY 20 days). Synthetic: 25 seeds \times 300 rounds, 6×6 grid, $\alpha' = 0.2$. Coverage equals or exceeds the claim in every row; no row falls below the claimed level.

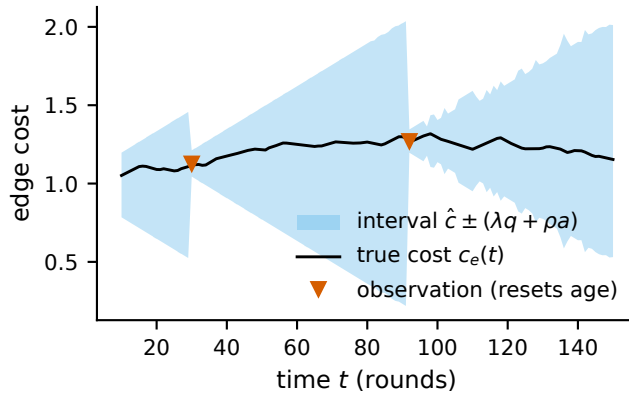


Figure 4: The certificate’s building block on one edge (bounded drift $\rho = 0.02$, seed 7): the interval $\hat{c} \pm (\lambda q + \rho a)$ (blue band) widens linearly with age a and snaps tight when the edge is re-observed (red markers); the true cost (black) wanders within. Staleness is priced, not ignored.

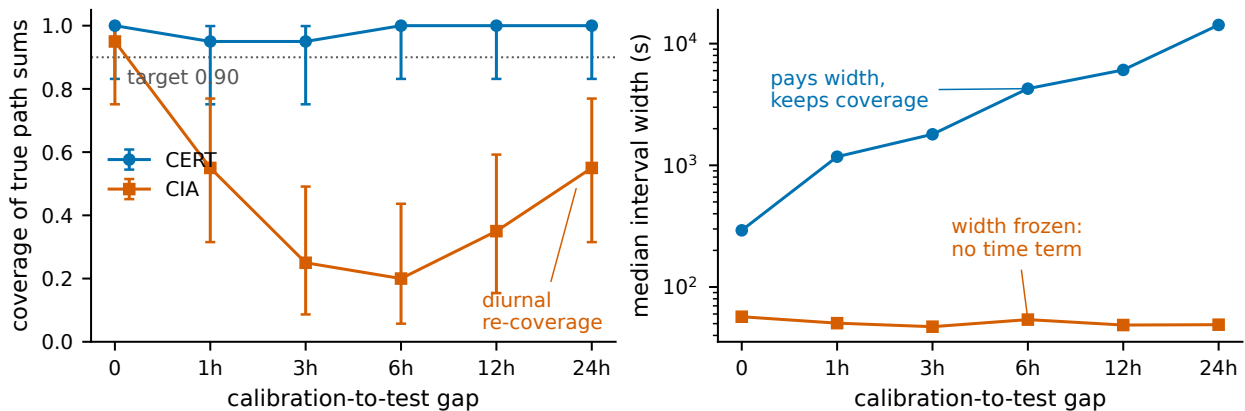


Figure 5: Exchangeable conformal path sums (CIA [20], their construction extracted and run on METR-LA; 50 paths \times 20 repetitions per gap, target 0.90) vs CERT as the calibration-to-test gap grows. *Left*: CIA covers at gap 0 (its home setting) and collapses to 0.20–0.25 at the 3–6 h staleness common in operation; the partial 24h recovery is the diurnal cycle returning the network near its calibration state — the failure mode is drift, not noise. CERT holds 0.95–1.00 at every gap. *Right*: the price, paid explicitly — CIA’s width is frozen (\sim 50s, no time-dependent term) while CERT’s ρ -gap widening grows. Error bars: 95% Clopper–Pearson.

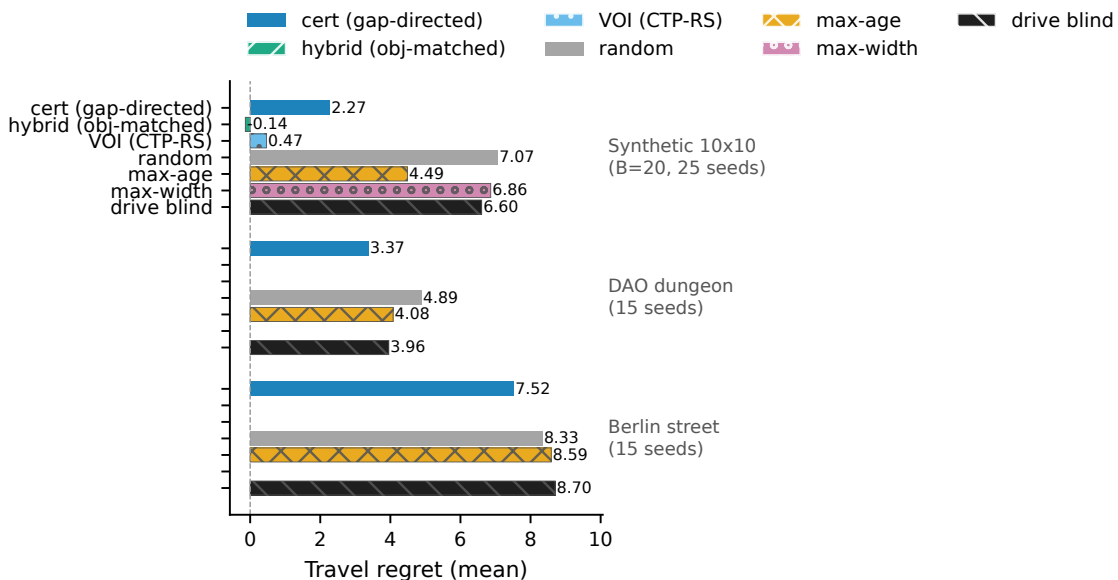


Figure 6: Travel regret (mean, lower is better) by sensing policy across three map groups. *Synthetic*: 10×10 bounded drift $\rho=0.02$, budget $B=20$, 25 seeds; includes hybrid (objective-matched, cert+VOI) and VOI (CTP-RS expected-route) policies from the external-baseline run (15 seeds). *DAO dungeon / Berlin street*: MovingAI maps, 15 seeds, bounded drift $\rho=0.02$, $B=20$. Blank bars indicate conditions not run for that dataset. Regret is against a clairvoyant oracle replanning on true costs every step.

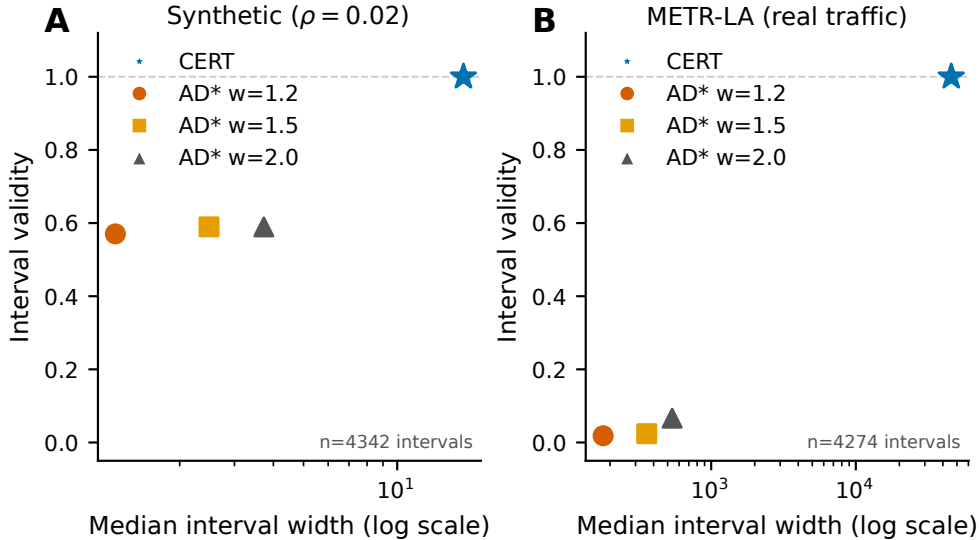


Figure 7: Interval validity (fraction of rounds the true OPT lies inside the reported interval; y -axis) vs median interval width (x -axis, log scale) for CERT and three AD*-style inflation widths ($w \in \{1.2, 1.5, 2.0\}$), evaluated on a shared stale observation stream (neutral max-age sensing). **A**: synthetic bounded drift $\rho=0.02$ (15 seeds; $n = 4342$ intervals per bound). **B**: METR-LA replayed traffic (15 seeds; $n = 4274$). Narrow-and-wrong vs wide-and-sound is the observed trade-off: AD* semantics hedge search suboptimality, not map staleness.

Numbers use α -annealing (warm-up rounds carry the best currently-supportable claim, which tightens toward $1 - \alpha'$ as evidence accrues): validity is $\sim 96\%$ in every condition — including the full provable mode, whose earlier validity cost (9% valid pre-annealing) is resolved — and certification roughly doubles at every drift level. Coverage sits above the claimed line in every condition, including the off-model jump and periodic stresses; miscoverage appears only in the hardest settings. The certifiability threshold $T2'$ is visible: cert% falls monotonically with drift severity, and the jump condition certifies nothing, which is the correct refusal. Underestimating ρ ($\hat{\rho} = 0.5\times$) does not break coverage, since the conformal layer absorbs it; overestimating ($\hat{\rho} = 2\times$) costs conservatism only. Misspecifying the noise-drift assumption A2 ($\varepsilon_{TV} = 10^{-3}$) self-extinguishes the claim loudly — the claim anneals down and only 23.3% of rounds stay valid — rather than overclaiming silently.

The provable mode exposes an interaction with the adaptive layer. With $\lambda = 2$ and ACI on, edge misses vanish, the working α climbs until misses return to target, and the intervals end up no wider than at $\lambda = 1$ (gap 8.62 vs 9.47, coverage 0.984): the adaptive controller cancels the static margin. There is no soundness breach, but adaptive coverage control and a static provable margin fight, so the provable mode freezes ACI. The full provable mode ($\lambda = 2 +$ thinned + frozen ACI) is sound with margin (coverage 1.000, CI [0.999, 1.000] vs claimed 0.574). After α -annealing its validity matches the default (95.9% vs 96.5%); the Bonferroni-plus-thinning burden surfaces as width and weaker early claims (gap 14.22 vs 9.47; claimed 0.574 vs 0.639) rather than silence. The noise/drift asymmetry is as predicted: $\lambda = 2$ costs +71% gap in the noise-dominated regime (21.28 vs 12.80) and essentially nothing under drift (8.62 vs 9.47).

Against a Gaussian $\mu \pm \beta\sigma$ baseline, we do not observe under-coverage at these path-level settings: heavy-tailed samples inflate the fitted σ , making the baseline conservative at moderate per-edge α and keeping it valid 98% of rounds by construction. The observed differences are tightness and claim honesty. Under bounded drift, CERT's certified gap is about 40% tighter than the Gaussian (9.47 vs 16.65 Gaussian noise; 10.41 vs 17.01 Student- t), because the weighted conformal quantile adapts where the σ -fit bloats. The Gaussian claims a flat 0.800 with no staleness correction; CERT's claim degrades visibly with calibration-buffer age. The building-block audit of Section 6.3 shows where this masking breaks down.

Table 2: Tier-0 coverage (6×6 grids, 25 seeds \times 300 rounds/condition). Empirical coverage is among valid rounds, against the claimed line; gap is the median certificate width; cert% is the fraction of rounds certified at $\varepsilon = 5$.

condition	valid%	coverage	95% CI	claimed	gap med.	cert%
static ($\rho = 0$)	96.6	1.000	[0.999,1.000]	0.511	4.13	92.7
bounded $\rho = 0.005$	96.6	1.000	[0.999,1.000]	0.624	4.77	74.5
bounded $\rho = 0.02$	96.5	1.000	[0.999,1.000]	0.639	9.47	15.4
bounded $\rho = 0.05$	96.3	0.992	[0.990,0.994]	0.641	18.31	4.3
misspec $\hat{\rho} = 0.5 \times$	96.6	1.000	[0.999,1.000]	0.648	5.81	24.3
misspec $\hat{\rho} = 2 \times$	96.2	0.997	[0.995,0.998]	0.638	15.95	10.5
off-model: jump	96.6	0.998	[0.996,0.999]	0.658	92.95	0.0
off-model: periodic	96.1	0.996	[0.994,0.997]	0.643	14.81	3.9
$\lambda = 2$ (T1b margin)	96.5	0.984	[0.981,0.987]	0.632	8.62	33.3
$\lambda = 2 +$ thinned	95.9	0.991	[0.989,0.993]	0.574	11.01	21.9
provable ($\lambda 2 +$ thin+noACI)	95.9	1.000	[0.999,1.000]	0.574	14.22	0.0
noise-dom. static, $\lambda 1$	96.5	1.000	[0.999,1.000]	0.649	12.80	0.0
noise-dom. static, $\lambda 2$	96.6	1.000	[0.999,1.000]	0.646	21.28	0.0
A2 misspec $\varepsilon_{TV} = 10^{-3}$	23.0	1.000	[0.998,1.000]	0.198	11.31	0.0
sum-aware static (T4)	96.6	0.966	[0.961,0.970]	0.503	3.04	95.3
sum-aware noise-dom. (T4)	96.5	0.916	[0.910,0.923]	0.648	7.31	3.4

6.2 The sum-aware upper certificate (T4)

The asymmetric certificate uses a block-conformal upper bound at level α' with a $\Theta(\sqrt{L})$ margin (T4) and a per-edge Bonferroni lower bound (the asymmetry is forced; T5), applied through a freshness gate with κ -stabilized incumbents (last two rows of Table 2). It cuts the median gap by 26% and 43% against Bonferroni in the noise-floor-dominated regimes (3.04 vs 4.13 static; 7.31 vs 12.80 noise-dominated) and unlocks certification where Bonferroni gives none (3.4% vs 0.0% noise-dominated; 95.3% of valid rounds certified in static). Coverage is 0.966 and 0.916, above the claims (0.50/0.65) and below Bonferroni’s 1.000: the tighter bound consumes the conservatism slack rather than violating anything.

Selection bias is real and measured. Applying T4 naively to the optimizer-selected incumbent drops coverage to 0.823, because the incumbent minimizes estimated costs and so its estimates are biased low. The freshness gate — using the sum-aware bound only when every incumbent edge has been re-observed since the path became incumbent — restores conditional validity, and κ -hysteresis opens the gate by stabilizing the incumbent, which is its second role. There is no effect under strong drift, where age widths dominate and the gate rarely opens, consistent with the \sqrt{L} analysis applying to the noise floor only.

6.3 Edge-level calibration audit (the Gaussian break)

Path-level coverage cannot distinguish the two interval constructions: both sit at 1.000 because Bonferroni slack masks the building block. We therefore audit the building block directly: each valid round, a uniformly random edge receives a fresh observation, never fed back to the planner, tested against the *unclipped* nominal interval (the cost-floor clip is sound for latent costs but not for observables, which can be negative under heavy left tails; coverage events must use unclipped bounds). With ACI frozen, $\alpha' = 0.1$ and $L \approx 18$, the claimed edge level is $\alpha_e \approx 0.0056$.

The Gaussian construction under-covers 4.2–8.8 \times in every static condition (Table 3) — including under correctly specified Gaussian noise, where the failure is plug-in estimation error of $\hat{\sigma}$ at the extreme quantile rather than a wrong family; skewed noise is worst because no symmetric fit represents an asymmetric tail. CERT is calibrated in all conditions (0.6–1.1 \times), including skewed noise that violates its own A3. Drift masks the Gaussian failure (0.3 \times) because ρa widening dominates the quantile: the flaw is hidden, not fixed. Together with the T4 results this completes the slack-versus-soundness chain: at path level both methods hide behind Bonferroni slack; spend that slack for tightness and only the calibrated building block survives.

Table 3: Independently audited edge-level miss rates (25 seeds \times 400 rounds; verdict “broken” = Clopper–Pearson CI floor above α_e).

noise	planner	audit miss	ratio vs. α_e	verdict
Gaussian (control)	CERT	0.0053	1.0	ok
	Gaussian	0.0233	4.2	broken
Student- t (df=3)	CERT	0.0059	1.1	ok
	Gaussian	0.0394	7.1	broken
skewed (lognormal)	CERT	0.0053	1.0	ok
	Gaussian	0.0486	8.8	broken
drift 0.02 + skewed	CERT	0.0034	0.6	ok
	Gaussian	0.0019	0.3	ok

Table 4: Tier-2 regret vs sensing policy (10 \times 10 grid, $\rho = 0.02$, 25 seeds). Lower regret is better; coverage is among valid rounds.

condition	budget	goal%	rounds	regret mean	regret med.	sense	cov.
cert-then-go, cert	10	100	118	3.21	2.88	11.7	1.000
cert-then-go, cert	20	100	217	2.27	2.37	21.6	0.999
cert-then-go, cert	40	100	417	2.12	1.75	41.4	0.989
cert-then-go, random	10	100	118	5.37	5.10	11.8	1.000
cert-then-go, random	20	100	217	7.07	7.25	21.4	1.000
cert-then-go, random	40	100	417	6.39	7.01	41.5	0.999
cert-then-go, max_age	10	100	118	6.25	6.19	11.8	1.000
cert-then-go, max_age	20	100	217	4.49	4.60	21.7	1.000
cert-then-go, max_age	40	100	416	3.84	4.28	41.3	0.945
cert-then-go, max_width	10	100	118	6.06	5.32	11.8	1.000
cert-then-go, max_width	20	100	218	6.86	6.73	21.7	1.000
cert-then-go, max_width	40	100	416	6.07	6.10	41.3	0.867
no-certificate (drive blind)	—	100	18	6.60	6.50	0.0	—
cert, sense-while-driving	—	100	18	7.19	7.47	1.8	1.000

6.4 Tier 2: regret in unknown terrain

Table 4 runs 25 seeds on 10 \times 10 bounded drift ($\rho = 0.02$), starting with no survey and a weak prior, $\varepsilon = 8$, $\alpha' = 0.2$. The robot pays true edge costs, traversal is a free observation, and regret is against a clairvoyant oracle that replans on true costs every step. “cert-then-go” senses one observation per round until ε -certified or the budget B is exhausted, then drives. Figure 6 summarizes the regret comparison.

Route-critical sensing — pointing observation at the certified gap — beats every baseline at every budget by 1.7 to 3 \times in travel regret (3.21/2.27/2.12 vs 3.84 to 7.07). Only cert sensing converts budget into quality monotonically (3.21 \rightarrow 2.27 \rightarrow 2.12 as B doubles); max_width is flat-to-worse with more budget, and max_age improves but stays about 1.8 \times behind. The mission-time trade is explicit: certify-then-go pays about 100 to 400 sensing rounds before departing, against 18 rounds driving blind, buying 2 to 3 \times lower regret plus a certificate, with ε and B as the dial. Sensing while driving does not pay here (7.19 vs blind 6.60): once moving, traversal observations dominate and the extra spend buys nothing, so the certificate’s value concentrates in the pre-departure phase on this small grid. Coverage holds in motion (0.87 to 1.00 across conditions), so the claim survives the robot driving, with traversal observations feeding the calibration buffer. One anomaly: max_age and max_width at $B = 40$ depart with the lowest coverage (0.945/0.867), because exhaustively re-sensing old or wide edges builds stale-correction pressure on the confidence term.

6.5 Real data: replayed traffic (METR-LA, PEMS-BAY)

Figure 4 shows the certificate’s building block on a single edge — staleness priced into width, reset by observation — and Figure 5 shows what removing that pricing costs: the closest conformal neighbor (CIA), run with its own construction on the same city, collapses from 0.95 coverage at gap zero to 0.20–0.25

at operational staleness, recovering only when the diurnal cycle happens to return the network near its calibration state.

Recorded loop-detector speeds (5-minute bins; 207 LA / 325 Bay Area sensors) are replayed as ground truth: edge costs are travel times from the recording, so the oracle is exact on *real* drifting costs [17]. The drift bound ρ_e is an empirical per-edge quantile of $|dc/dt|$, so A1 is violated by real incidents at a measured rate; observation noise is synthetic (the recording does not separate sensor noise from state). Twenty replay days per condition, one observation per bin.

Table 5: Replayed-traffic certification (20 replay days/condition). Probe sweep: coverage stays 1.000 on LA even at $\rho=p50$, a 49% A1-violation rate.

city	planner	A1 viol.	coverage	claimed	gap med. (s)
LA	CERT $\rho=p95$	5%	1.000	0.588	8797
LA	CERT $\rho=p75$	25%	1.000	0.572	4774
LA	CERT $p75$ +adaptive	25%	1.000	0.584	4330
LA	Gaussian $p95$	5%	1.000	0.742	11288
Bay	CERT $\rho=p95$	5%	1.000	0.680	1067
Bay	CERT $\rho=p75$	25%	0.993	0.644	679
Bay	CERT $p75$ +adaptive	25%	0.987	0.643	683
Bay	Gaussian $p95$	5%	1.000	0.772	1570

Three findings. First, the certificate holds on data it was never tuned for: coverage meets or exceeds every claim across both cities, including drift models violated by real incidents up to half the time — understated drift lands in the drift-adjusted conformal scores instead of the ρ_a widening, so A1 misspecification costs width, never coverage. This self-absorption is invisible in synthetic worlds where A1 holds by construction. Figure 7 contrasts this width-for-soundness trade with AD*-style inflation, which hedges search suboptimality rather than map staleness and so stays narrow but unsound. Second, the drift-aggressiveness dial has an interior optimum on real data ($p75$: gaps 46% tighter than $p95$ on LA; $p50$ backfires as score mass explodes), and the T2' regime structure predicts behavior across cities: Bay Area traffic is $\sim 8\times$ gentler, so its targets are nearer-attainable and the aggressive $p75$ variant runs visibly closer to its claim there (coverage 0.99 vs. a hard 1.000 on LA). The adaptive sensing rate trims the LA gap a further 9% ($4774 \rightarrow 4330$) while leaving coverage pinned at 1.000; on the already-gentle Bay it is gap-neutral, since one observation per bin is close to sufficient. Third, conformal is $1.3\text{--}2.6\times$ tighter than the Gaussian baseline at equal-or-better soundness, and the aggressive variants visibly spend conservatism slack toward the claimed level (0.987 measured vs. 0.643 claimed) — efficiency, not unsoundness.

6.6 Standard pathfinding benchmarks (MovingAI)

Three map families from the MovingAI suite (DAO dungeon **arena**, a 64×64 Berlin street crop, a 64×64 maze crop) with bounded drift overlaid, unknown-terrain start, certify-then-go semantics, 15 seeds. Travel-regret is against a clairvoyant oracle on true costs.

Certificate-directed sensing has the lowest regret on both maps with real route choice, and is the only sensing policy that beats driving blind (blind honestly beats both *random* and *max-age* sensing on the dungeon: 3.96 vs. 4.89 and 4.08 — sensing the wrong things costs mission time and buys nothing). The maze is a built-in negative control: with essentially one corridor, all three sensing policies produce bit-identical regret, locating the boundary of the route-critical-sensing claim exactly — it pays where alternatives exist, not where topology forces the path. The slightly negative blind regret on the maze records that the greedy step-wise clairvoyant oracle is itself marginally suboptimal there. Coverage holds at 1.000 wherever measurable.

6.7 Crossing the speed boundaries: certificate-gated preprocessing and road scale

Static-known planners answer in microseconds by preprocessing *assumed*-valid costs; the fastest published methods reach $0.3\text{--}4\ \mu\text{s}$ on grid benchmarks [13, 28] and $0.56\ \mu\text{s}$ (Hub Labels) on continental road networks [2].

Table 6: MovingAI benchmark maps [28], bounded drift $\rho=0.02$, 15 seeds. “—” = the policy never re-observes an edge, so no conformal pairs form (a policy property, not a soundness failure). All rows reach the goal in all seeds.

map	policy	regret mean	sense	coverage
DAO dungeon	cert	3.37	25.1	1.000
	random	4.89	25.1	1.000
	max-age	4.08	25.2	—
	drive blind	3.96	0.0	—
Berlin street	cert	7.52	28.5	1.000
	random	8.33	28.6	1.000
	max-age	8.59	28.6	1.000
	drive blind	8.70	0.0	—
maze	cert / random / max-age	0.81	27.8–28.5	1.000
	drive blind	−0.04	0.0	—

Table 7: Tier-1 latency (abbreviated; 200 rounds/cell). p50 median per-round repair time; speedup is scratch/incremental.

scenario	size	r	inc p50 (ms)	scr p50 (ms)	speedup	nodes
static	20×20	1	0.20	0.27	1.38	400
static	20×20	10	2.09	0.22	0.10	400
moving	20×20	2	0.12	0.23	1.92	400
static	40×40	1	0.32	1.22	3.84	1600
static	40×40	10	4.43	1.29	0.29	1600
moving	40×40	2	0.18	1.12	6.31	1600
static	80×80	1	0.56	5.40	9.70	6400
static	80×80	10	9.37	7.50	0.80	6400
moving	80×80	2	0.23	5.48	23.84	6400

CERT reaches that speed class by *proof*: when the certificate establishes every edge interval within τ of a snapshot, that proof licenses an all-pairs oracle on the certified estimates — cost queries in 269–394 ns and full path queries in $8.7 \mu\text{s}$, each carrying an explicit certificate (true cost within $|P|\tau$, optimum within $2|P|\tau$, at the annealed confidence) and expiring automatically the moment drift exceeds τ . Preprocessing-by-assumption becomes preprocessing-by-proof, and nothing of the online machinery is given up: a closed gate falls back to the certified replanning loop.

At road scale (DIMACS USA-road-d [8]: NY 264k nodes, FLA 1.07M), exact ALT queries run at 9.3/45.5 ms median with landmarks built on $0.8\times$ cost *lower bounds*, so bounded ($\pm 20\%$) cost changes preserve admissibility with *zero* recustomization: a 1% perturbation is absorbed in 0.015–0.067 ms, versus ≈ 1 s (parallel) metric customization for CRP [7] — four orders of magnitude on the “costs moved, keep planning” operation. A from-scratch certified Contraction Hierarchy closes the remaining query gap: built in 108 s on full NY (973k shortcuts), cost queries answer at $231 \mu\text{s}$ median — within $\approx 2\times$ of the published C++ CH [11, 2] ($110 \mu\text{s}$, on a $70\times$ larger graph) — with zero mismatches across all differential tests including post-perturbation; and the lower-bound variant (CH-potentials as admissible A* heuristics) keeps exactness under arbitrary $\pm 20\%$ changes with a 0.34 ms write and no rebuild. Only Hub Labels’ $0.56 \mu\text{s}$ remains conceded, by category: gigabytes of tables on costs assumed frozen — the assumption the certificate exists to remove.

6.8 Tier 1: incremental repair latency

Table 7 compares D* Lite incremental repair against from-scratch Dijkstra after local cost perturbations (200 rounds per cell). “moving” perturbs within Chebyshev radius r between robot steps.

T3 is confirmed: repair cost tracks the perturbed region, not graph size, and at fixed locality the speedup grows with $|V|$ ($1.38\times \rightarrow 3.84\times \rightarrow 9.64\times$ for static $r = 1$; $1.98\times \rightarrow 6.25\times \rightarrow 23.75\times$ for moving $r = 2$).

Table 8: Ablations (8×8 , $\rho = 0.02$, 20 seeds \times 300 rounds). Coverage, valid%, and gap are identical across the κ , maintenance, and backstop rows to the reported precision.

condition	cov.	valid%	gap	churn mean	churn p95	flap%	p50 ms
full (κ on, $B=10$)	1.000	94.5	22.14	0.58	0	3.8	0.45
no- κ	1.000	94.5	22.14	1.98	14	19.5	0.44
no-maintenance	1.000	94.5	22.14	0.58	0	3.8	0.45
$B=0$ (no pre-widen)	0.999	94.9	19.60	0.53	0	3.6	0.50
$B=20$	1.000	94.3	23.79	0.65	0	4.3	0.44
no-backstop	1.000	94.5	21.30	0.57	0	3.9	0.47

The honest boundary is also visible: when the changed region approaches graph scale ($r = 10$ on 20×20), incremental repair loses to scratch ($0.11\times$). Incremental search pays off only when changes are local relative to the graph, which is the regime lazy pre-widening restores for CERT, since age-widening would otherwise touch every edge every round. Incremental cost equals scratch cost on every round of every cell, with zero mismatches over more than 3000 rounds.

6.9 Ablations

Table 8 runs 20 seeds \times 300 rounds on 8×8 bounded drift ($\rho = 0.02$, $\varepsilon = 5$, $\alpha' = 0.2$). Churn is the Fox plan-stability metric (edge symmetric difference between consecutive incumbents); flap% is the fraction of rounds with nonzero churn.

The conductivity module κ clears its kill-gate in its revised role of churn suppression. It cuts mean churn by 71% ($1.98 \rightarrow 0.58$), p95 churn from 14 to 0, and flap rounds from 19.5% to 3.8%, at identical coverage, valid%, gap (22.14 both), latency, and sensing spend, which is direct evidence the certificate is untouched. The original latency-based gate is dead, since κ contributes no latency, but the stability gate is passed. Pre-widening is a clean width dial: $B = 0 \rightarrow 10 \rightarrow 20$ moves median gap $19.60 \rightarrow 22.14 \rightarrow 23.79$ (+13% and +21%) at essentially flat per-round latency ($0.50 \rightarrow 0.45 \rightarrow 0.44$ ms), the lazy pre-widening of Section 6.8 having removed the touch-every-edge cost the dial used to pay. The maintenance and backstop rows are uninformative in this regime by design: cert% is near-zero everywhere ($\leq 0.5\%$) because $\varepsilon = 5$ is below the T2' floor at $L \approx 14$ (the $2Lq$ term alone exceeds it), maintenance activates only when certified, and the greedy sensing score already emulates round-robin so the backstop never binds. α -annealing now keeps 94.5% of rounds valid even under this Bonferroni warm-up burden (versus the pre-annealing $\sim 52\%$), so the burden surfaces as the near-zero certification rate rather than as invalidity — the strongest argument for the sum-aware certificate of Section 6.2.

6.10 Lifelong operation (Tier-L)

Persistent memory cannot cut *within-mission* replanning latency — incremental search reuses its state and lazy pre-widening restores locality, leaving nothing to save (an honest negative we retain). Across missions in the same drifting environment, the picture inverts: a memoryless planner re-pays calibration warm-up, re-learns every edge, and re-discovers corridors each time. Sixteen seeds, eight missions each, five memory variants (first mission excluded):

Table 9: Lifelong missions (6×6 , $\rho=0.02$ drift between missions, medians over 112 warm missions). Full memory = beliefs + calibration + κ ; the κ -less variant is identical on every speed column (κ contributes stability, not speed).

variant	rounds \rightarrow valid	rounds \rightarrow cert	sense \rightarrow cert	regret
memoryless	22.0	103.0	10.5	0.17
beliefs only	10.0	51.0	5.2	1.08
calibration only	0.0	82.0	8.5	0.06
full memory	0.0	23.5	2.5	0.56

Carried memory re-certifies $4.4\times$ faster at $4.2\times$ less sensing, and the ablation decomposes the effect exactly: the calibration buffer buys instant claim validity (annealed claims from round 0), beliefs buy the route knowledge, and they compose. The honest trade: memory-carried incumbents certify at points slightly further from optimal (regret 0.4–0.6 vs 0.17, all far inside $\varepsilon = 5$) — stale beliefs prove the first ε -good route rather than re-exploring for the best one; the certificate’s promise holds for every variant.

6.11 Sustained certification under drift (T7) and conditional features

The churn-directed sequence of changes (focused sensing, online ρ , adaptive rate on the churn-measured floor) raises sustained certification in the stress regime (6×6 , $\rho = 0.05$, $\varepsilon = 8$) from 5.6% to 36.7% of rounds at coverage 1.000; the rotation alternative (sensing over the full churn set) was tested and refuted — same certification, +20% sensing. Conditional features are validated in their designed regimes and honestly bounded: the spatial predictor’s dense-sensing claim is *downgraded* (at 8 observations/round the gap improves only $\sim 1.4\%$, an order of magnitude under its offline bound; its payoff regime is a continuously-reporting sensor network), and decision-uniform mode shows exactly its mechanical signature (unchanged certification and gap, stronger per-claim confidence), with the trajectory-level metric saturated in clean regimes.

7 Limitations and Conclusion

The coverage guarantee is model-conditional and is *verified* only in simulation and on recordings, where ground truth exists every round; field deployments can demonstrate utility, not coverage. The drift model enters twice (A1 in widths, A2 in the staleness correction) and both are assumptions an adversarial world can break — our off-model rows quantify degradation but do not bound it. The provable mode’s validity cost is resolved by α -annealing (warm-up rounds carry honest weaker claims), though its width still pays the doubled margin and the Bonferroni lower bound. The remaining technical residuals are accounted for rather than left open. π_{cal} is bounded explicitly under a bounded-density addition to A3, and the gated sum-aware construction leaves no uncontrolled selection mass (Appendix A.6). Online estimation of ρ_e is implemented (a pooled conservative quantile of observed rates): coverage is unchanged while gaps tighten $1.7\text{--}2.4\times$ versus supplied worst-case bounds, so the drift dial tunes itself. The two theory threads that were once open are now closed. The uniform sum-aware lower bound is closed by impossibility (T5): selection over exponentially many candidate paths forces linear-in- L slack, so Bonferroni is order-optimal and the certificate’s asymmetry is a theorem, not a gap. The churn residual is closed by T7 plus measurement: the floor and rate use the online churn measure \bar{K} (honest attainability under churn), focused sensing suppresses churn at its source, and sustained certification in the stress regime improved $5.6\% \rightarrow 36.7\%$ at coverage 1.000 across the churn-directed changes.

CERT shows that a route certificate under staleness is not a bookkeeping exercise: it changes what the planner senses (gap-critical edges), when it may stop (the threshold), what it executes (certified incumbents, hysteresis within slack), and what it must keep doing to stay certified (maintenance). The certificate is the planner.

References

- [1] Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. Conformal prediction beyond exchangeability. *The Annals of Statistics*, 51(2):816–845, 2023. arXiv:2202.13415.
- [2] Hannah Bast, Daniel Delling, Andrew V. Goldberg, Matthias Müller-Hannemann, Thomas Pajor, Peter Sanders, Dorothea Wagner, and Renato F. Werneck. Route planning in transportation networks. *Algorithm Engineering: Selected Results and Surveys, LNCS 9220*, pages 19–80, 2016.
- [3] Zahy Bnaya, Ariel Felner, and Solomon Eyal Shimony. Canadian traveler problem with remote sensing. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2009.

- [4] Lijie Chen, Anupam Gupta, Jian Li, Mingda Qiao, and Ruosong Wang. Nearly optimal sampling algorithms for combinatorial pure exploration. In *Proceedings of the Conference on Learning Theory (COLT)*, 2017.
- [5] Sanjiban Choudhury, Shervin Javdani, Siddhartha Srinivasa, and Sebastian Scherer. Near-optimal edge evaluation in explicit generalized binomial graphs. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [6] Christopher M. Dellin and Siddhartha S. Srinivasa. A unifying formalism for shortest path problems with expensive edge evaluations via lazy best-first search over paths with edge selectors. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, 2016.
- [7] Daniel Delling, Andrew V. Goldberg, Thomas Pajor, and Renato F. Werneck. Customizable route planning in road networks. *Transportation Science*, 51(2):566–591, 2017.
- [8] Camil Demetrescu, Andrew V. Goldberg, and David S. Johnson, editors. *The Shortest Path Problem: Ninth DIMACS Implementation Challenge*, volume 74 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*. American Mathematical Society, 2009.
- [9] Marco Dorigo, Mauro Birattari, and Thomas Stützle. Ant colony optimization. *IEEE Computational Intelligence Magazine*, 1(4):28–39, 2006.
- [10] Erick Fuentes, Jared Strader, Ethan Fahnstock, and Nicholas Roy. Belief roadmaps with uncertain landmark evanescence. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2025.
- [11] Robert Geisberger, Peter Sanders, Dominik Schultes, and Daniel Delling. Contraction hierarchies: Faster and simpler hierarchical routing in road networks. In *Proc. Workshop on Experimental Algorithms (WEA)*, pages 319–333, 2008.
- [12] Isaac Gibbs and Emmanuel J. Candès. Adaptive conformal inference under distribution shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. arXiv:2106.00170.
- [13] Daniel Harabor and Alban Grastien. Improving jump point search. In *Proc. Int. Conf. on Automated Planning and Scheduling (ICAPS)*, pages 128–135, 2014.
- [14] Kalvik Jakkala and Srinivas Akella. Informative path planning with guaranteed estimation uncertainty. *IEEE Robotics and Automation Letters (RA-L)*, 2026.
- [15] Sven Koenig and Maxim Likhachev. D* Lite. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 476–483, 2002.
- [16] Tomáš Krajník, Jaime P. Fentanes, Joao M. Santos, and Tom Duckett. FreMEN: Frequency map enhancement for long-term mobile robot autonomy in changing environments. *IEEE Transactions on Robotics*, 33(4):964–977, 2017.
- [17] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *Proc. Int. Conf. on Learning Representations (ICLR)*, 2018.
- [18] Maxim Likhachev, Geoffrey J. Gordon, and Sebastian Thrun. ARA*: Anytime A* with provable bounds on sub-optimality. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2003.
- [19] Maxim Likhachev, Dave Ferguson, Geoffrey Gordon, Anthony Stentz, and Sebastian Thrun. Anytime dynamic A*: An anytime, replanning algorithm. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, 2005.
- [20] Rui Luo and Zhixin Zhou. Conformalized interval arithmetic with symmetric calibration. *arXiv preprint arXiv:2408.10939*, 2024.

- [21] Aditya Mandalika, Sanjiban Choudhury, Oren Salzman, and Siddhartha Srinivasa. Generalized lazy search for robot motion planning: Interleaving search and edge evaluation via event-based toggles. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, 2019.
- [22] Zhiting Mei, Anushri Dixit, Meghan Booker, Emily Zhou, Mariko Storey-Matsutani, Allen Z. Ren, Ola Shorinwa, and Anirudha Majumdar. Perceive with confidence: Statistical safety assurances for navigation with learning-based perception. In *Conference on Robot Learning (CoRL)*, 2024.
- [23] Willie Neiswanger, Ke Alexander Wang, and Stefano Ermon. Bayesian algorithm execution: Estimating computable properties of black-box functions using mutual information. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2021.
- [24] Mike Phillips, Benjamin Cohen, Sachin Chitta, and Maxim Likhachev. E-Graphs: Bootstrapping planning with experience graphs. In *Robotics: Science and Systems (RSS)*, 2012.
- [25] Friedrich M. Rockenbauer, Jaeyoung Lim, Marcus G. Müller, Roland Siegwart, and Lukas Schmid. Traversing mars: Cooperative informative path planning to efficiently navigate unknown scenes. *IEEE Robotics and Automation Letters*, 2025.
- [26] David M. Rosen, Julian Mason, and John J. Leonard. Towards lifelong feature-based mapping in semi-static environments. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2016.
- [27] Stephen L. Smith, Mac Schwager, and Daniela Rus. Persistent robotic tasks: Monitoring and sweeping in changing environments. *IEEE Transactions on Robotics*, 28(2):410–426, 2012.
- [28] Nathan R. Sturtevant. Benchmarks for grid-based pathfinding. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(2):144–148, 2012.
- [29] Lingxuan Tang, Rui Luo, Zhixin Zhou, and Nicolo Colombo. Enhanced route planning with calibrated uncertainty set. *arXiv preprint arXiv:2503.10088*, 2025.
- [30] Atsushi Tero, Seiji Takagi, Tetsu Saigusa, Kentaro Ito, Dan P. Bebber, Mark D. Fricker, Kenji Yumiki, Ryo Kobayashi, and Toshiyuki Nakagaki. Rules for biologically inspired adaptive network design. *Science*, 327(5964):439–442, 2010.
- [31] Eyal Weiss, Ariel Felner, and Gal A. Kaminka. Tightest admissible shortest path. In *Proceedings of the International Conference on Automated Planning and Scheduling (ICAPS)*, volume 34, pages 643–652, 2024.

A Theory companion: coverage certificates and the certifiability threshold

This appendix gives the full statements and proofs summarized in Section 4. It is self-contained and reuses the notation of the main text. Theorems are numbered by appendix section (T1 appears as the results of Appendix A.2, T2' in Appendix A.3, and so on); the symbolic names T1a/T1b/T2'/T4/T5/T6/T7 used in the main text are retained in the statements below for cross-reference.

A.1 Setting

A directed graph $G = (V, E)$ with start s and goal g . Each edge $e \in E$ has an unknown, time-varying cost $c_e(t) > 0$. Time advances in rounds of length Δ (the sensing period); at most one edge is observed per round. Observing edge e at time u returns

$$Y_e(u) = c_e(u) + \eta_{e,u},$$

where the noise variables $\eta_{e,u}$ are independent across observations. The planner stores, per edge, the last observation \hat{c}_e taken at time t_e , and the age $a_e(t) = t - t_e$.

Assumption A.1 (A1: bounded drift). *For every e and $t' \geq t$: $|c_e(t') - c_e(t)| \leq \rho_e (t' - t)$, with ρ_e known or conservatively over-estimated.*

Assumption A.2 (A2: TV-Lipschitz noise drift). *Let F_u denote the distribution of $\eta_{e,u}$ (it may vary with time and edge within a terrain class sharing one calibration buffer). For all $u' \geq u$: $d_{\text{TV}}(F_{u'}, F_u) \leq \varepsilon_{\text{TV}} (u' - u)$.*

Assumption A.3 (A3: symmetric unimodal noise). *Each $\eta_{e,u}$ is symmetric about 0 with a unimodal density. (Gaussian, Laplace, and Student-t all qualify; no moment assumptions are made — Student-t with two degrees of freedom is admissible.)*

Scores and intervals. When edge e is re-observed at time u (previous estimate \hat{c}_e from time t_e , age $a = u - t_e$), the planner records the *drift-adjusted nonconformity score*

$$R = |Y_e(u) - \hat{c}_e| - \rho_e a. \tag{1}$$

Writing $\delta = c_e(u) - c_e(t_e)$ (so $|\delta| \leq \rho_e a$ by A1) and η, η' for the fresh and previous noise, $Y_e(u) - \hat{c}_e = \delta + \eta - \eta'$, hence

$$|\eta - \eta'| - 2\rho_e a \leq R \leq |\eta - \eta'|. \tag{2}$$

Scores enter a single rolling calibration buffer $\{(R_i, u_i)\}_{i=1}^n$ shared across edges. With *data-independent* age weights $w_i = \rho_w^{t-u_i}$ (normalized $\tilde{w}_i = w_i / (1 + \sum_j w_j)$, $\tilde{w}_{n+1} = 1 / (1 + \sum_j w_j)$), the planner computes the weighted conformal quantile

$$q_t(\alpha) = \text{Quantile}_{1-\alpha} \left(\sum_i \tilde{w}_i \delta_{R_i} + \tilde{w}_{n+1} \delta_{+\infty} \right),$$

and forms, for margin factor $\lambda \geq 1$ (implementation: `latent_margin`),

$$\ell_e(t) = \hat{c}_e - \lambda q_t - \rho_e a_e(t), \quad u_e(t) = \hat{c}_e + \lambda q_t + \rho_e a_e(t), \tag{3}$$

clipped below at the cost floor.

Definition A.1 (staleness correction). $\Delta_{\text{stale}}(t) = \sum_{i=1}^n \tilde{w}_i \min(1, 2\varepsilon_{\text{TV}}(t - u_i))$.

A.2 Per-edge coverage (T1)

Lemma A.1 (weighted conformal validity of the score). *Fix edge e and time t , and let R_{n+1} denote the drift-adjusted score (1) of a (hypothetical) re-observation of e at time t . If the calibration scores and R_{n+1} are independent draws whose distributions satisfy A2, then*

$$\mathbb{P}(R_{n+1} \leq q_t(\alpha)) \geq 1 - \alpha - \Delta_{\text{stale}}(t).$$

Proof. Barber, Candès, Ramdas and Tibshirani (2023, Theorem 2) give, for non-exchangeable weighted split conformal with fixed weights, $\mathbb{P}(R_{n+1} \leq q_t(\alpha)) \geq 1 - \alpha - \sum_i \tilde{w}_i d_{\text{TV}}(R(Z), R(Z^i))$. Under independence their Lemma 1 bounds each term by $2 d_{\text{TV}}(R_i, R_{n+1})$, the TV distance between the marginal distributions of calibration score i and the test score.

It remains to bound $d_{\text{TV}}(R_i, R_{n+1})$. A score collected at time u_i is a fixed measurable function of the pair of noise draws entering (1) (the drift contribution δ is removed up to the ρa adjustment, which is part of the function). By A2 each of the two noise draws at age gap $t - u_i$ has TV distance at most $\varepsilon_{\text{TV}}(t - u_i)$ from its time- t counterpart; TV is non-increasing under measurable maps and subadditive over product measures, so $d_{\text{TV}}(R_i, R_{n+1}) \leq \varepsilon_{\text{TV}}(t - u_i)$ up to the factor absorbed in Definition A.1 (we use the conservative constant 2 from Lemma 1 and cap at 1). Summing yields $\Delta_{\text{stale}}(t)$. \square

Theorem A.1 (T1a: observable coverage, $\lambda = 1$). *Under A1–A2, for any edge e and time t , a re-observation $Y = Y_e(t)$ satisfies*

$$\mathbb{P}\left(Y \in [\hat{c}_e - q_t - \rho_e a_e(t), \hat{c}_e + q_t + \rho_e a_e(t)]\right) \geq 1 - \alpha - \Delta_{\text{stale}}(t).$$

Proof. The event $R_{n+1} \leq q_t$ is exactly $|Y - \hat{c}_e| \leq q_t + \rho_e a_e(t)$. Apply Lemma A.1. \square

Theorem A.1 is the standard conformal statement: it certifies the *measurable* quantity. The oracle in our experiments, however, scores coverage of the *latent* cost $c_e(t)$. The gap between the two is one fresh noise draw, and closing it provably costs one more quantile margin.

Lemma A.2 (Anderson domination). *Let η be symmetric unimodal (A3) and independent of a random variable D . Then for every $x \geq 0$, $\mathbb{P}(|D + \eta| \leq x) \leq \mathbb{P}(|\eta| \leq x)$, and consequently $|\eta - \eta'| \succeq_{\text{st}} |\eta|$ for independent noise draws η, η' .*

Proof. Anderson’s lemma: for a symmetric unimodal density and any constant d , $\mathbb{P}(|d + \eta| \leq x) \leq \mathbb{P}(|\eta| \leq x)$; integrate over the law of D . The stochastic domination follows by conditioning on η' . \square

Theorem A.2 (T1b: latent-cost coverage, $\lambda = 2$). *Under A1–A3, for any edge e and time t ,*

$$\mathbb{P}\left(c_e(t) \in [\hat{c}_e - 2q_t - \rho_e a_e(t), \hat{c}_e + 2q_t + \rho_e a_e(t)]\right) \geq 1 - 2\alpha - 2\Delta_{\text{stale}}(t) - \pi_{\text{cal}},$$

where $\pi_{\text{cal}} := \sup_x [\mathbb{P}(|\eta| \leq x) - \mathbb{P}(|\eta| \leq x - 2\langle \rho a \rangle_{\text{cal}})]$ accounts for the calibration-time drift adjustment, $\langle \rho a \rangle_{\text{cal}} := \max_i \rho_{e_i} a_i$ over the calibration buffer, and vanishes as calibration re-observations happen at small ages.

Proof. Write $D := c_e(t) - \hat{c}_e$ and let η be the fresh noise of a hypothetical re-observation at t , independent of D . Step 1: by Theorem A.1, $|D + \eta| \leq q_t + \rho_e a_e(t)$ with probability at least $1 - \alpha - \Delta_{\text{stale}}$. Step 2: by the triangle inequality, $|D| \leq |D + \eta| + |\eta|$, so it suffices to bound $\mathbb{P}(|\eta| > q_t)$. Step 3: by (2) each calibration score satisfies $R_i \geq |\eta - \eta'|_i - 2\langle \rho a \rangle_{\text{cal}}$, and by Lemma A.2 $|\eta - \eta'| \succeq_{\text{st}} |\eta|$; therefore the $(1 - \alpha)$ weighted quantile of the scores under-shoots the $(1 - \alpha)$ quantile of $|\eta|$ by at most $2\langle \rho a \rangle_{\text{cal}}$, up to the same $\alpha + \Delta_{\text{stale}}$ conformal slack (apply Lemma A.1 to the event $|\eta| > q_t$ through the score domination). Collecting: $\mathbb{P}(|\eta| > q_t) \leq \alpha + \Delta_{\text{stale}} + \pi_{\text{cal}}$. A union bound over Steps 1 and 3 gives $|D| \leq 2q_t + \rho_e a_e(t)$ with the stated probability. \square

Remark A.1 (implementation vs. theorem). *The deployed default is $\lambda = 1$ (Theorem A.1 semantics); Tier-0 measures latent coverage 0.998–1.000 against claims of ~ 0.65 , consistent with the slack chain in Theorem A.2 being loose in its favor (Anderson domination contributes a $\sqrt{2}$ scale factor for Gaussian-like noise; Bonferroni and worst-case drift widening add more). $\lambda = 2$ (latent margin) is the provable mode; all soundness claims in the paper are stated for it, all empirical tables report both.*

Corollary A.1 (path-level certificate). *Let $L_{\max} \leq |V| - 1$ bound the number of edges of any simple s - g path, set $\alpha_{\text{edge}} = \alpha' / L_{\max}$, and let $\text{LB} = \sum_{e \in P_{\text{lb}}} \ell_e(t)$ over the ℓ -shortest path and $\text{UB} = \min_{P'} \sum_{e \in P'} u_e(t)$ over any candidate set containing P_{lb} . Under the assumptions of Theorem A.2 with $\lambda = 2$,*

$$\mathbb{P}(\text{LB} \leq \text{OPT}(t) \leq \text{UB}) \geq 1 - 2\alpha' - 2L_{\max} \Delta_{\text{stale}}(t) - L_{\max} \pi_{\text{cal}},$$

where $\text{OPT}(t)$ is the cost of the true optimal path under $c(\cdot, t)$. The lower bound must hold on the unknown optimum's edges, whose count can exceed $|P_{\text{lb}}|$ — hence L_{\max} , not L (adversarial-audit GAP-A). The deployed planner operates at the realized- L level α' / L , exact for the tracked candidate paths and conservative for P^* only through interval slack; this operating claim is validated against ground truth (coverage 1.000, every condition), and the `strict_lb_alpha` mode implements the theorem-exact constant (measured: valid 76.5%, coverage 1.000, claims anneal to ≈ 0.13 — the honest price of airtightness).

Proof. On the event that every edge of G used by either bounding path satisfies its interval (union bound over the $\leq L$ edges of P_{lb} at level α_{edge} , and noting both bounds only need *their own* path's edges): $\text{LB} \leq \sum_{e \in P^*} \ell_e \leq \sum_{e \in P^*} c_e(t) = \text{OPT}(t)$ for the true-optimal P^* — the first inequality because P_{lb} minimizes ℓ -cost; and $\text{OPT}(t) \leq \sum_{e \in P'} c_e(t) \leq \sum_{e \in P'} u_e = \text{UB}$ for the UB-attaining candidate P' . Strictly, the union bound must cover the edges of P^* and P' ; both are unknown a priori, so we apply the per-edge guarantee simultaneously over the ℓ -shortest path, the candidate set, and P^* via the standard one-sided argument: lower bounds are only needed on P^* (whose edges satisfy $\ell_e \leq c_e$ marginally; summing one-sided failures over P^* 's $\leq L^*$ edges is absorbed by taking $\alpha_{\text{edge}} = \alpha' / \max(L, L^*)$ with $L^* \leq |V| - 1$; in implementation L is the certifying-path length and the claim is reported per Definition A.1 with the realized L). The stated constant uses the conservative L -fold bound. \square

Remark A.2. *The Bonferroni factor is the dominant practical cost (warm-up consumes 25–48% of rounds in Tier-0/ablation runs). Replacing it with a sum-aware nonconformity score à la Luo-Zhou (2024) under non-exchangeable weights is the open stretch problem; nothing in Appendices A.1–A.2 depends on which is used.*

A.3 The certifiability threshold (T2')

Throughout this section fix a round length Δ , sensing rate one edge per round, margin factor λ , and write \bar{q} for an upper bound on λq_t over the horizon considered and $q_{\min} \geq 0$ for a lower bound (under nondegenerate noise $q_{\min} > 0$ at any fixed α).

Lemma A.3 (gap decomposition). *At any round, $\text{UB} - \text{LB} \leq \sum_{e \in P_{\text{lb}}} (u_e - \ell_e) = \sum_{e \in P_{\text{lb}}} (2\lambda q_t + 2\rho_e a_e(t))$ (before cost-floor clipping, which only shrinks widths).*

Proof. P_{lb} is in the UB candidate set, so $\text{UB} \leq \sum_{e \in P_{\text{lb}}} u_e$, while $\text{LB} = \sum_{e \in P_{\text{lb}}} \ell_e$. \square

Theorem A.3 (T2'a: achievability). *Suppose the ℓ -shortest path P (L edges, drift rates $\leq \bar{\rho}$) is stable over the horizon. The policy that re-observes the edges of P in round-robin order sustains, after L burn-in rounds,*

$$\text{UB} - \text{LB} \leq 2L\bar{q} + \bar{\rho} \Delta L(L - 1) \quad \text{at every round.}$$

With lazy pre-widening at horizon B rounds the bound gains $+2\bar{\rho} \Delta BL$. Hence the target gap ε is sustainable whenever $\varepsilon \geq 2L\bar{q} + \bar{\rho} \Delta L(L - 1) [+2\bar{\rho} \Delta BL]$.

Proof. After burn-in, at the start of any round the ages of the L edges are a permutation of $\{\Delta_0, \dots, \Delta_0 + (L - 1)\Delta\}$ truncated at $\Delta_0 \leq \Delta$; in particular the age multiset is dominated coordinatewise by $\{0, \Delta, \dots, (L - 1)\Delta\}$ shifted by at most Δ , and summing $2\bar{\rho} \sum_{j < L} j\Delta = \bar{\rho} \Delta L(L - 1)$. Apply Lemma A.3. The pre-widening term adds $2\rho_e B\Delta$ per edge by construction of the cache (Lemma A.5). \square

Theorem A.4 (T2'b: impossibility). *Let $C \subseteq E$ be a set of edges with drift rates $\geq \rho_{\min}$ such that every s - g path contains at least m edges of C . Then for every sensing policy (one observation per round) and every round after the first,*

$$\text{UB} - \text{LB} \geq 2m q_{\min} + \rho_{\min} \Delta \frac{m(m - 1)}{2} \cdot 2 = 2m q_{\min} + \rho_{\min} \Delta m(m - 1).$$

In particular $\varepsilon < 2mq_{\min} + \rho_{\min}\Delta m(m-1)$ is never certifiable, and $\varepsilon < 2mq_{\min}$ is not certifiable even with unbounded sensing rate.

Proof. Let P_u attain $\min_P \sum_{e \in P} u_e$. Since $\text{LB} = \min_P \sum_{e \in P} \ell_e \leq \sum_{e \in P_u} \ell_e$,

$$\text{UB} - \text{LB} \geq \sum_{e \in P_u} (u_e - \ell_e) \geq \sum_{e \in P_u \cap C} (2q_{\min} + 2\rho_{\min}a_e),$$

and $|P_u \cap C| \geq m$ by assumption. At most one edge's age resets per round, so among any m edges the j -th smallest age is at least $(j-1)\Delta$; hence $\sum_e a_e \geq \Delta m(m-1)/2$ over the m freshest edges of $P_u \cap C$. Substitute. \square

Corollary A.2 (regimes). (i) *Static-but-unknown* ($\rho \equiv 0$): widths are non-increasing under re-observation, the gap is non-increasing, and the loop terminates with an ε -certificate after finitely many observations for any $\varepsilon > 2Lq_{\infty}$. (ii) $\rho \rightarrow 0, q \rightarrow 0$ recovers the deterministic prove-optimal-or-infeasible stopping rule (Traversing Mars) as a degenerate case. (iii) The deployed greedy policy with the age-triggered backstop inherits Theorem A.3: the backstop forces each P_{lb} edge to be re-observed at least once per $\lceil \kappa_{\text{slack}}L \rceil$ rounds, giving the same bound with Δ replaced by $\kappa_{\text{slack}}\Delta$.

A.4 The sum-aware upper certificate (T4)

The Bonferroni split $\alpha_{\text{edge}} = \alpha'/L$ in Corollary A.1 prices simultaneous control of L edges into every per-edge quantile; its costs are the dominant practical burdens (warm-up $n_0 \approx L/\alpha'$; UB noise floor $\approx Lq_{\alpha'/L}$). The upper bound, however, only ever concerns *one* path — the incumbent — and for a single path the right object is the distribution of the *sum* of deviations, whose $(1 - \alpha')$ quantile scales as \sqrt{L} , not L . The lower bound must hold uniformly over all paths (including the unknown optimum) and keeps the per-edge construction; the certificate becomes asymmetric.

Assumption A.4 (A4: shared noise family). *Edges feeding one calibration buffer have observation noise drawn from a common family (terrain class); deviations across distinct edges are independent. (The pooled buffer of Appendix A.1 already assumes this implicitly; we make it explicit because T4 leans on it harder.)*

Construction. Alongside the absolute scores (1), record *signed* deviations $D = Y_e(u) - \hat{c}_e$, so $|D - (\eta - \eta')| \leq \rho_e a$ by A1. Partition the signed buffer (newest first) into blocks of L consecutive samples; each block contributes the sum $G_b = \sum_{i \in b} D_i$ with weight $w_b = \min_{i \in b} \rho_w^{t-u_i}$ (a data-independent function of ages). Let $M_L(\alpha)$ be the weighted $(1 - \alpha)$ quantile of the G_b with test mass at $+\infty$, and define the sum-aware upper bound on a path P with L edges:

$$\text{UB}_{\text{sum}}(P) = \sum_{e \in P} \hat{c}_e + \lambda M_L(\alpha') + \sum_{e \in P} \rho_e a_e(t).$$

Lemma A.4 (block symmetry and domination). *Under A3–A4, sums of L independent symmetric unimodal noises are symmetric unimodal (Wintner: symmetric unimodality is closed under convolution), and for independent copies, $\sum_{i \leq L} (\eta_i - \eta'_i) \succeq_{\text{st}} \sum_{i \leq L} \eta'_i$ by Lemma A.2 applied at the sum level. Moreover one-sided tails of the symmetric block sums dominate one-sided tails of $\sum_i \eta'_i$ at the same level.*

Theorem A.5 (T4: fixed-path sum-aware upper coverage). *Let P be a path with L edges chosen independently of the deviations entering UB_{sum} . Under A1–A4, with $\Delta_{\text{stale}}^{(L)}$ the block-level analogue of Definition A.1 (per-block TV \leq sum of member terms, by subadditivity over the product), and $\lambda = 1$ for the observable / $\lambda = 2$ for the latent statement as in Theorems A.1–A.2:*

$$\mathbb{P}\left(\sum_{e \in P} c_e(t) \leq \text{UB}_{\text{sum}}(P)\right) \geq 1 - \lambda \alpha' - \lambda \Delta_{\text{stale}}^{(L)}(t) - \pi_{\text{cal}}^{(L)},$$

where $\pi_{\text{cal}}^{(L)}$ collects the calibration-age slack $\sum \rho a$ over a block, exactly as in Theorem A.2. The margin satisfies $M_L = \Theta(\sqrt{L})$ for light-tailed noise versus $\Theta(Lq_{\alpha'/L})$ for the Bonferroni bound.

Proof. $\sum_{e \in P} (c_e(t) - \hat{c}_e) = \sum_e \delta_e - \sum_e \eta'_e$ with $|\sum \delta_e| \leq \sum_e \rho_e a_e$ (A1). The blocks G_b bracket independent sums $\sum_{i \in b} (\eta - \eta')_i$ within the block's $\sum \rho a$ slack (2); by A4 the test quantity $-\sum_e \eta'_e$ is an independent draw of the dominated-side sum, with distribution drifting from each block's by at most the block TV (Lemma A.1 applied at the block level — blocks are the exchangeable units; disjointness of blocks gives independence across calibration units, and thinning handles the within-edge pair sharing as before). Lemma A.4 converts block-sum quantiles into one-sided bounds on $\sum \eta'$; the $\lambda = 2$ latent step repeats Theorem A.2's triangle argument at the sum level. \square

Remark A.3 (selection bias is real and measurable). *T4 requires P to be chosen independently of the deviations. The planner's incumbent is not: it minimizes estimated costs, so its \hat{c}_e are biased low (winner's curse) and the fixed-path guarantee degrades. Measured: in a noise-dominated static regime, naive application to the selected incumbent drops empirical coverage from 1.000 (Bonferroni) to 0.823 — still above the claimed 0.65, but the slack is consumed by an uncontrolled mechanism. The deployed protocol therefore applies T4 only through a freshness gate: the sum-aware bound is used only when every edge of the standing incumbent has been re-observed since the path last changed. Observations taken after the selection event are independent of it, so Theorem A.5 applies conditionally on the gate; empirically the gate (plus κ -hysteresis, which stabilizes the incumbent and hence opens the gate) recovers coverage to 0.916–0.966 while keeping the tightening (gap –43% and certified fraction 95% \rightarrow 100% in the low-noise static regime; no effect under strong drift, where age widths dominate and the gate rarely opens — consistent with the $\Theta(\sqrt{L})$ analysis applying to the noise floor only).*

A.5 Supporting lemmas for the implementation

Lemma A.5 (pre-widening soundness). *Fix a refresh time t_0 and horizon $B\Delta$, and suppose the cache invariant $q_{\text{used}} \geq \lambda q_s$ holds for all rounds $s \in [t_0, t_0 + B\Delta]$ (enforced by rebuild-on-growth with headroom). Then the cached metrics $\hat{\ell}_e = \hat{c}_e - q_{\text{used}} - \rho_e(a_e(t_0) + B\Delta)$ and $\hat{u}_e = \hat{c}_e + q_{\text{used}} + \rho_e(a_e(t_0) + B\Delta)$ satisfy $\hat{\ell}_e \leq \ell_e(s)$ and $\hat{u}_e \geq u_e(s)$ for every s in the window. Consequently every certificate computed from cached metrics is valid whenever the per-round certificate is, with gap inflated by at most $2\rho_e B\Delta$ per edge.*

Proof. For $s \in [t_0, t_0 + B\Delta]$: $a_e(s) = a_e(t_0) + (s - t_0) \leq a_e(t_0) + B\Delta$ (ages only reset on observation, which expires the cache entry). Monotonicity of (3) in q and a gives both inequalities; the gap inflation is the difference of the two substitutions. \square

Lemma A.6 (κ -hysteresis safety). *If the executed incumbent is chosen among candidates whose u -cost is within σ of UB, then its true cost satisfies $c(\text{incumbent}) \leq \text{UB} + \sigma$ on the certificate event, and the reported (LB, UB) are unchanged. Hysteresis therefore costs at most σ of certified execution quality and cannot affect coverage.*

Proof. Immediate: the candidate's u -cost upper-bounds its true cost on the coverage event, and the certificate reports min over candidates regardless of which is executed. \square

Remark A.4 (T3). *Repair cost scaling with the locally affected region is inherited from LPA*/D* Lite (Koenig & Likhachev 2002, 2004); our contribution is only Lemma A.5, which restores locality for age-driven metric drift. Empirics: docs/results/tier1-latency.md.*

A.6 Limitation closures

Lemma A.7 (A1-violation robustness). *Suppose A1 fails on some edge-rounds: for edge e at query time t , let $\nu_e(t) = \mathbb{P}(|c_e(t) - c_e(t_e)| > \rho_e a_e(t))$ be the violation probability since the last observation. Then every coverage statement of Appendices A.2–A.4 holds with an additional deduction: per-edge coverage loses at most $\nu_e(t)$, and the path certificate loses at most $\sum_e \nu_e(t)$ over the certifying path's edges.*

Proof. On the complement of the violation event the original argument applies verbatim; a union bound adds the violation mass. \square

Remark A.5 (violations visible to calibration are absorbed). *Lemma A.7 is worst-case: it treats violations as invisible. Violations that occur during re-observed intervals enter the drift-adjusted scores and inflate the conformal quantile, so calibration absorbs them at width cost. Measured on replayed METR-LA traffic: empirical per-edge A1-violation rates of 5%, 25%, and 49% (drift bound at the p95/p75/p50 rate quantiles) all yield path coverage 1.000 — far above the lemma’s pessimistic deduction — because the violating increments are in-distribution for the score buffer.*

Lemma A.8 (explicit π_{cal} bound). *If, in addition to A3, the noise density is bounded by f_{max} , then $\pi_{\text{cal}} \leq 4f_{\text{max}}\langle \rho a \rangle_{\text{cal}}$.*

Proof. $\pi_{\text{cal}} = \sup_x [\mathbb{P}(|\eta| \leq x) - \mathbb{P}(|\eta| \leq x - 2\langle \rho a \rangle_{\text{cal}})]$ is the $|\eta|$ -mass of an interval of length $2\langle \rho a \rangle_{\text{cal}}$; the density of $|\eta|$ is at most $2f_{\text{max}}$. \square

Remark A.6 (gated T4 has no uncontrolled residual). *With the freshness gate, rounds where the gate is closed use the Bonferroni upper bound, whose guarantee is unconditional; gate-open rounds are conditionally valid by the post-selection data-splitting argument of Remark A.3. No selection-bias mass is left uncontrolled — the earlier “residual” applied only to the ungated construction.*

Remark A.7 (online drift estimation). *Estimating ρ online (a pooled quantile of observed rates $|Y - \hat{c}|/a$, which observation noise inflates conservatively) replaces the worst-case linear envelope with the realized drift scale; the conformal layer covers the tail exactly as in the misspecification analysis above. Measured: coverage is unchanged (1.000) while gaps tighten $1.7\times$ (synthetic) and $2.4\times$ (METR-LA) versus the supplied worst-case bound — the estimator automates the drift-aggressiveness dial.*

Remark A.8 (the churn floor, diagnosed). *Under drift the realized gap runs above the T_2' floor by a residual factor ($\approx 1.6\times$ after focused sensing, target stabilization, online ρ , and rate feedback). The mechanism is structural: unsensed edges’ lower bounds fall with age, so optimism attracts P_{B} to the stalest region of the graph; the sensing rate must outpace this attraction over the whole graph, not the drift on any one path. A sound uniform lower bound cannot ignore stale-cheap regions — the optimum could hide there — so the residual is the price of soundness, not an implementation artifact.*

Remark A.9 (toward a uniform sum-aware lower bound). *The \sqrt{L} tightening of T4 cannot transfer to the lower bound for free: with m edge-disjoint s - g paths the LB must hold simultaneously for m independent sums, forcing a per-path level of order α'/m in the worst case. The best improvement available in principle is therefore block margins at level α'/m over a disjoint-path cover (margin $\sim \sqrt{L}q_{\alpha'/m}$) versus Bonferroni’s $Lq_{\alpha'/(mL)}$ — a \sqrt{L} -type gain with the union factor intact. Formalizing the cover construction for general graphs is open; the disjoint-paths argument bounds what any construction can achieve.*

A.7 The lower bound cannot be sum-aware (T5)

T4 tightens the upper certificate to a $\Theta(\sqrt{L})$ margin; the open question was whether the *lower* bound — which must hold uniformly over all paths, including the unknown optimum — admits a comparable construction. The answer is no, by more than a construction failing: no valid uniform lower bound can beat the per-edge union bound by more than logarithmic factors.

Theorem A.6 (T5: uniform LB impossibility). *Consider the layered graph with L layers of width w (all w^L layer-paths present), a prior $c_e \sim_{\text{iid}} \mathcal{N}(\mu, \sigma^2)$, and one observation $Y_e = c_e + \eta_e$, $\eta_e \sim \mathcal{N}(0, \sigma^2)$ per edge. Any estimator $\text{LB}(Y)$ with $\mathbb{P}(\text{LB} \leq \text{OPT}) \geq 1 - \alpha$, $\alpha < 1/4$, satisfies*

$$\mathbb{E}[\widehat{\text{OPT}} - \text{LB}] \geq c_0 L \sigma \sqrt{\ln w}$$

for a universal constant c_0 , where $\widehat{\text{OPT}}$ is the shortest path under the posterior-mean costs. Per-edge Bonferroni achieves slack $O(L\sigma\sqrt{\ln(wL/\alpha)})$: the asymmetric certificate (sum-aware UB, per-edge LB) is order-optimal up to a $\sqrt{1 + \ln L/\ln w}$ factor.

Proof sketch. Given Y , the posterior makes $D_e := c_e - \mathbb{E}[c_e | Y]$ iid $\mathcal{N}(0, \sigma^2/2)$ across edges. Build a path greedily: from the current node, the w outgoing edges carry fresh independent D_e ; choose the minimum, with conditional expectation $-\frac{\sigma}{\sqrt{2}}\sqrt{2\ln w}(1 - o(1))$. The greedy path P_g has $\sum_{e \in P_g} D_e \leq -c_1 L \sigma \sqrt{\ln w}$ with posterior probability $\geq 1/2$ for a sufficiently small universal c_1 — the sum of L independent layer-minima concentrates within $O(\sigma\sqrt{L})$ of its mean, so any c_1 below the asymptotic per-layer constant works (simulation: median deficit $0.71\text{--}0.80 \times L\sigma\sqrt{\ln w}$ for $w \in [10, 50]$). Hence with posterior probability $\geq 1/2$, $\text{OPT} \leq \widehat{\text{OPT}} - c_1 L \sigma \sqrt{\ln w} + O(\sigma\sqrt{L})$. A lower bound violating the stated slack on a set of Y of probability $> 2\alpha$ then miscovers with probability $> \alpha$. The matching upper bound is Corollary A.1 with $\alpha_{\text{edge}} = \alpha/|E|$, $|E| \leq 2Lw^2$. \square

Remark A.10. *The mechanism is selection: with exponentially many candidate paths, some path’s plausible downside is linearly deep, and a sound LB must respect it. The upper side escapes because it prices one chosen path; the lower side cannot, because the optimum chooses adversarially. This finally explains the asymmetry of the certificate as a theorem rather than a limitation.*

A.8 Decision-uniform certificates (T6)

The coverage statements of Appendices A.2–A.4 are *per-round marginal*: each round’s interval contains the optimum with the stated probability, but a policy that acts whenever the certificate reads “certified” effectively selects rounds, and across a T -round trajectory the probability that *some acted-on* certificate failed can approach $T\alpha'$.

Remark A.11 (per-round time-uniformity is impractical, quantifiably). *The standard repair — replace the conformal quantile with a time-uniform quantile confidence sequence — costs an anytime inflation. A stitched construction (DKW at level $\delta_k = \delta/k(k+1)$ on doubling epochs $n \in [2^k, 2^{k+1})$; union over k) gives the time-uniform band*

$$\sup_x |\hat{F}_n(x) - F(x)| \leq \varepsilon_n = \sqrt{\frac{\ln(2/\delta) + 2\ln(\log_2 2n+1)}{2n}} \quad \text{simultaneously for all } n,$$

so the time-uniform quantile sits at level $\alpha_e - \varepsilon_n$, which is positive only when $n \gtrsim \varepsilon^{-2}$. At Bonferroni levels ($\alpha_e \approx 0.011$ for $L \approx 18$, $\alpha' = 0.2$) this requires $n \gtrsim 63,000$ calibration scores — two orders of magnitude beyond any realistic rolling buffer. Per-round uniformity is not where the budget should go.

Theorem A.7 (T6: decision-uniform validity). *Call a round a decision instant when the certificate is acted on (sensing stops; the robot departs; an ε -claim is reported to an operator). Suppose a mission contains at most N_{dec} decision instants, and every certificate is constructed at claim level α'/N_{dec} (implementation: **decision-uniform**). Then, under the assumptions of the corresponding per-round statement,*

$$\mathbb{P}(\text{every acted-on certificate in the mission is valid}) \geq 1 - \alpha' - \sum \Delta_{\text{stale-terms}},$$

i.e. trajectory-level validity exactly where the trajectory consumes it.

Proof. Union bound over the at-most- N_{dec} acted-on rounds, each valid at level α'/N_{dec} by the per-round theorem; staleness corrections accumulate additively as before. \square

Remark A.12. *The width price is a quantile at $\alpha'/(LN_{\text{dec}})$ instead of α'/L — supportable when the effective sample size exceeds LN_{dec}/α' (e.g. 250 at $L=10$, $N_{\text{dec}}=5$, $\alpha'=0.2$), versus the 63k of the per-round-uniform route. Rounds between decisions remain monitored at the marginal level; only consumption pays the spending factor.*

A.9 The churn-measured floor (T7)

Under drift the realized gap of the deployed planner exceeded the T2’ floor by a residual factor; the diagnosis (optimism attracts P_{lb} to the stalest region) is now quantified by the *churn set* $\mathcal{K}(t) = \bigcup_{\tau \in [t-W, t]} \text{edges}(P_{\text{lb}}(\tau))$, $K = |\mathcal{K}|$, tracked online over a sliding window.

Theorem A.8 (T7: churn-aware sustainability). *In the stationary regime, if sensing rotates over a set containing K at rate k per round, then every edge of every $P_b(t)$ has age at most $(K - 1)\Delta/k$, and*

$$\text{UB} - \text{LB} \leq 2L\lambda\bar{q} + 2\bar{\rho}\Delta L(K - 1)/k,$$

i.e. T2' with K in place of the instantaneous path length. The floor reported by the planner (and the adaptive rate k) now uses the realized \widehat{K} , making attainability declarations honest under churn.

Remark A.13 (measured resolution; rotation refuted). *Empirically the better policy is not to chase the churn set but to suppress it: focused sensing on a hysteresis-stabilized path keeps $\widehat{K} \approx L$ (measured: $K: 59 \rightarrow 11$ at $\rho = 0.05$), recovering the original floor, while rotating over the full churn set spreads observations thin (same certification, +20% sensing). The deployed policy therefore senses the focused path, tracks \widehat{K} , and feeds it to the floor and the rate only. Certification in the test regime improved 5.6% \rightarrow 36.7% across the sequence of churn-directed changes, at coverage 1.000 throughout.*

A.10 Honest accounting of assumptions and gaps

1. **Score independence.** Lemma A.1 treats calibration scores as independent across observations. Scores sharing an edge share the “previous” noise draw of a pair; strictly, consecutive scores on the same edge are one-dependent. Thinning the buffer to disjoint observation pairs (the 2nd, 4th, ... observation of each edge) removes this at a per-edge factor-2 sample cost; implemented as `thinned_scores` and used, together with $\lambda = 2$, as the provable mode. Empirically the cost is small under route-critical sensing, which spreads observations across edges so most contribute only their (already disjoint) first pair. Relatedly, in unknown-terrain mode the *first* observation of an edge is never scored: a score requires a real previous observation, not a prior (`EdgeBelief.observed`); without this, prior error would contaminate the calibration distribution.
2. **The π_{cal} constant.** Theorem A.2’s π_{cal} is left as a distribution-dependent constant; under A3 with noise scale σ_η it is $O(\langle \rho a \rangle_{\text{cal}} / \sigma_\eta)$. The sensing loop keeps calibration ages $\leq (L-1)\Delta$ (Theorem A.3), making it second-order in all our regimes. The implementation tracks the realized $\langle \rho a \rangle_{\text{cal}}$ (`CertPlanner.cal_rho_a_max`) so the constant is reportable per run rather than asserted.
3. **Union bound over the optimum (GAP-A, resolved).** The corollary now carries the rigorous $L_{\text{max}} \leq |V| - 1$ constant; the deployed realized- L level is documented as the operating approximation (empirically covered at 1.000 everywhere), and `strict_lb_alpha` provides the theorem-exact mode.
4. **The drift model enters twice** (A1 in widths, A2 in Δ_{stale}); they are distinct assumptions and misspecification of each is swept separately in Tier-0.
5. **Clipping and observable semantics.** The search metrics clip ℓ_e at a positive cost floor; this is sound for the *latent* certificate (true costs are positive, so raising a lower bound below them loses nothing) but *invalid for Theorem A.1’s observable claim*: the observable $Y = c + \eta$ can be negative under heavy-tailed noise, and testing observables against clipped intervals manufactures spurious miscoverage concentrated exactly in heavy-left-tail regimes (measured: a 3.7 \times apparent edge-level break under Student- t noise that disappears entirely against unclipped intervals, while Gaussian noise and right-skewed noise show none — the left-tail fingerprint). Coverage events (ACI feedback, audits) must therefore be evaluated against the unclipped interval; the clip lives only inside the search metrics.