

Stage 11 / Step 1 — Go/NoGo Checklist

Warp → Detect → Denoise — readiness criteria for running Step 1 tests after establishing a dominant cognition well.

1) Preconditions

- 1 Calibration set collected for the target domain/task (100–5,000 prompts).
- 2 Tap layer selected (mid-late block, e.g., L-3..L-1).
- 3 PCA(3) whitener fitted on hidden states; funnel profile constructed (quantile fit + core deepening).

2) Measurements to capture (from the warped manifold)

- 1 Phantom Index (PI): separation between best and second best basin minima over the 2D PC plane.
- 2 Margin (Δ): $z(\min) - z(\min)$ in the warped surface; positive means the top basin is strictly deeper.
- 3 Well Depth Profile: normalized depth $\phi(r)$ and slope $g(r)$ over radius.
- 4 Radius Trace $r(t)$ & Well Score $S(t)$: from shadow decoding runs (no intervention).

3) Go/NoGo Thresholds (suggested)

Metric	Go (Proceed to Step 1)	Borderline (Tune)	NoGo (Rewarp)
Phantom Index (PI)	≤ 0.07 (target ≈ 0.06)	0.07–0.10	> 0.10
Margin Δ	≥ 0.04	0.02–0.04	< 0.02
Radius Trace $r(t)$	Monotonic \downarrow trend on shadow	Mostly \downarrow with small rebounds	Flat or \uparrow ; frequent rebounds
Well Score $S(t)$	Median ≥ 0.6 during reasoning spans	0.45–0.6	< 0.45

Notes:

- Thresholds reflect prior Stage 11 runs: $PI \approx 0.065$ and $\Delta \approx 0.044$ corresponded to stable single-well behavior and clean Step 1 outcomes.
- Use conservative 'Go' criteria for new domains; relax only after repeated stability.

4) Step 1 Test Procedure (once Go criteria met)

- 1 Shadow Run: enable Warp only; record $r(t)$, $S(t)$, matched-filter peaks (no rescoring).
- 2 Detect: run Stage 10 parser with null-calibrated dual thresholds; verify precision/recall vs baseline.
- 3 Denoise: enable EMA+median and phantom-guard; confirm $r(t)$ stabilizes and peaks align with reasoning spans.
- 4 Light-Touch Rescoring (optional in Step 1): $\alpha \leq 0.5$, $K \leq 16$; confidence-gate by $S(t)$ to avoid phantoms.

5) Step 1 Pass/Fail Gates (quick)

- 1 Precision ≥ 0.80 , Recall ≥ 0.98 on the chosen slice (or within 2–3% of prior Step 1 references).
- 2 Hallucination ≤ 0.26 and trending \downarrow versus baseline.

- 3 Abstain rate stable (no runaway abstention); phantom index does not increase post-denoise.

6) Troubleshooting if No-Go or Fail

- 1 Tap Scan: try layers $L-5..L-1$ and pick best PI/Δ pair.
- 2 Re-fit Funnel: increase core deepening or isotropize XY plane.
- 3 Tighten Nulls: raise absolute gate (q) or increase circular shift count K .
- 4 Back off Rescoring: reduce α or disable it entirely for Step 1.

Appendix — Run Record (fill per domain/model)

Model		Tap Layer		Date	
Calibration set (N)		PI		Margin Δ	
$S(t)$ median		$r_{\blacksquare}(t)$ trend		Rescoring α/K	
Precision		Recall		Hallucination	
Abstain rate		Outcome		Notes	