

# Chapter 4 — Hierarchical Bayesian Anchors for Long-Tail Estimation in Sparse Data

A production-grade prior system for pricing, time-to-sell, and tail risk under severe sparsity

Revision date: 2026-01-14

**Abstract.** This chapter formalizes the hierarchical anchor mechanism used in our marketplace models as a practical, leak-safe approximation to hierarchical Bayesian inference. The core idea is to compute, at each listing's time-zero  $t_0$ , a small set of robust, cohort-conditioned priors (anchors) for price and time-to-sell (TTS), and then express the listing's current state as relative deviation from those priors. Anchors stabilize learning in the two regimes that dominate high-frequency resale markets: (i) sparsity (most fine-grained cohorts have few recent SOLD observations) and (ii) heavy tails (a minority of listings persist for weeks, distorting naïve means). We show how the anchor system implements: (1) a backoff cascade across nested cohorts, (2) an embargo that prevents leakage from near-future outcomes, and (3) explicit support and level features that allow downstream models to condition on prior reliability.

## Key points.

- Why anchors exist: A direct model of absolute price or absolute duration is unstable under sparse cohorts; relative-to-anchor representations are stable and portable across generations and regions.
- What makes the anchors Bayesian: Each cohort estimate is shrunk toward broader cohorts according to support; the chosen backoff level is recorded as `anchor_level_k`.
- What makes the anchors leak-proof: SOLD comps are restricted to timestamps strictly before  $t_0$  (with an embargo), and anchor computation is executed in T0-certified, audited feature-store closures.
- Why this matters for the tail: Tail events (slow sales) are better separated when the model can compare the current ask against an appropriate, cohort-specific fair-value prior.

## 4.1 Problem setting: priors under high-frequency sparsity

Marketplace datasets look large in aggregate, yet are small where it matters: the moment we condition on the attributes buyers actually care about (generation, model tier, storage, condition, damage, location), the effective sample size collapses. At the same time, the observable we wish to predict—time-to-sale—has a long right tail. A small subset of listings remains active for weeks or months, and these “zombies” create high-variance estimates if we naïvely compute cohort averages.

We address this with a design principle: do not predict “absolute” quantities when a stable “relative” coordinate exists. In our system, the relative coordinate is a set of anchors computed at time-zero  $t_0$ : a robust estimate of the cohort’s recent fair price, and a robust estimate of the cohort’s typical time-to-sell. A listing’s price signal becomes the deviation between its ask and the anchor price; its liquidity signal becomes the deviation between its early trajectory and the anchor speed.

$$\text{PTV}_{\text{anchor}}(t_0) = \ln\left(\frac{\text{ask}(t_0)}{\text{anchor}(t_0)}\right)$$

Figure 4.1 — Price-to-value (PTV) as a log ratio between the ask price at  $t_0$  and an anchor (fair-value) price computed strictly from pre- $t_0$  SOLD comps.

## 4.2 A hierarchical Bayesian view of anchors

Anchors can be understood as empirical-Bayes posteriors for cohort-level latent parameters. Let  $y$  denote a robustly transformed outcome observed for SOLD comps—e.g., log price or log duration. Each cohort  $k$  has an unknown latent parameter  $\mu_k$  (e.g., a cohort’s typical log price), and cohorts are nested in a hierarchy where specific cohorts inherit information from broader ones.

A minimal generative model (for exposition) assumes: Within-cohort observations:  $y_i | \mu_k \sim \text{Normal}(\mu_k, \sigma^2)$  Hierarchy prior:  $\mu_k | \mu_{\text{parent}} \sim \text{Normal}(\mu_{\text{parent}}, \tau^2)$  The posterior mean for  $\mu_k$  is a convex combination of the cohort estimate and the parent prior. The weight on data grows with cohort support  $n_k$ ; under low support, the estimate shrinks to the parent.

$$\mu_k^{\text{post}} = w_k \bar{y}_k + (1 - w_k) \mu_{\text{parent}}, \quad w_k = \frac{n_k}{n_k + n_0}$$

Figure 4.2 — A canonical empirical-Bayes shrinkage rule. Our production anchors implement the same principle using robust medians/quantiles plus a discrete backoff cascade.

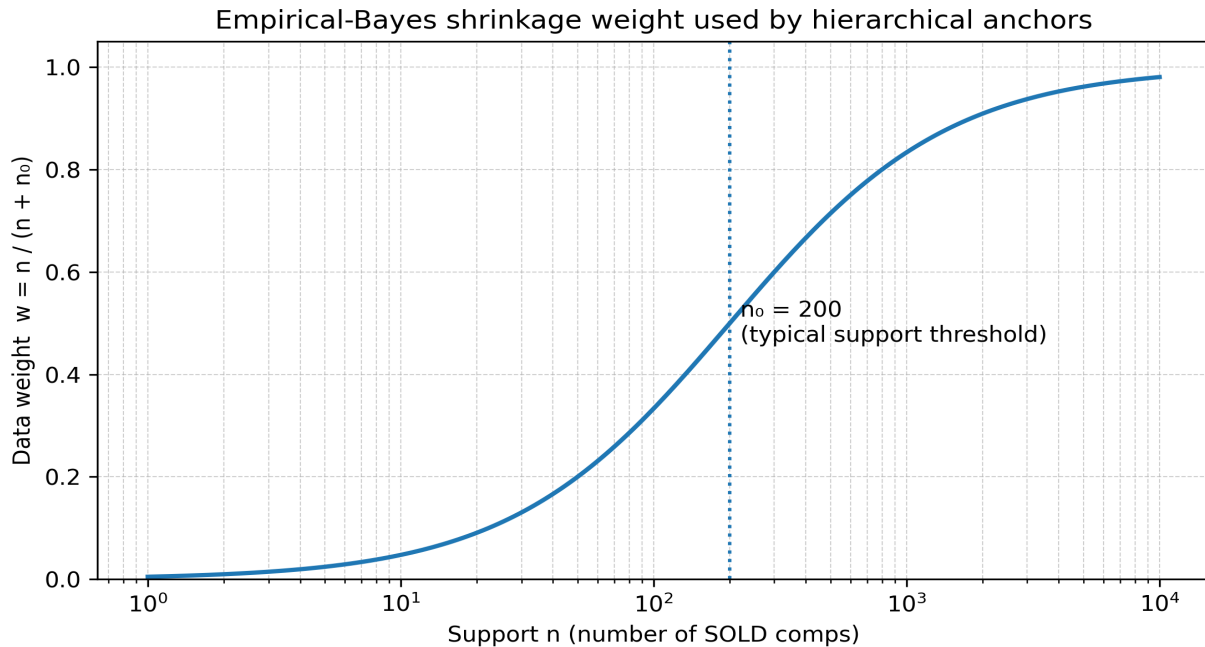
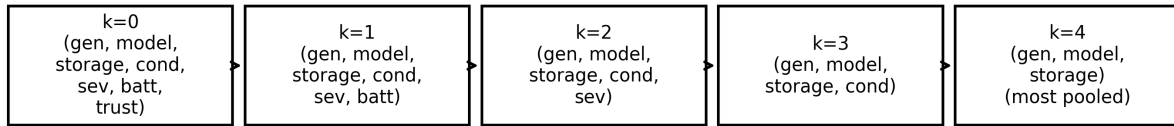


Figure 4.3 — Shrinkage weight as a function of cohort support  $n$ . The system behaves like a Bayesian posterior: small cohorts borrow strength from broader priors.

## 4.3 From theory to production: robust estimators and discrete backoff

Pure Bayesian inference is rarely the bottleneck; data engineering and leakage safety are. Our anchor implementation makes three deliberate choices: (1) robust statistics (medians / quantiles) instead of means, (2) a discrete backoff cascade instead of continuous hyperparameter inference, and (3) explicit, inspectable support metadata so downstream models can reason about uncertainty.

Advanced anchor (smart) backoff cascade  
 Select the most specific cohort with sufficient SOLD support; otherwise pool more aggressively



Hard locks: generation  $\times$  model\_norm  $\times$  storage\_bucket are never relaxed.  
 anchor\_level\_k  $\in \{0, \dots, 4\}$  records how much pooling was required.

Figure 4.4 — Backoff cascade used by the smart anchor. The anchor\_level\_k feature records how much pooling was required; hard locks (generation  $\times$  model\_norm  $\times$  storage bucket) are never relaxed.

## 4.4 Strict anchor at $t_0$ : recent-history priors with embargo

The strict anchor is designed for high precision when data is available. It uses recent SOLD history in 30-day and 60-day windows and computes robust medians for: (i) fair price, (ii) time-to-sell, and (iii) cohort support. Crucially, all SOLD comps used for these medians are constrained to occur strictly before  $t_0$ , with an embargo (e.g., 5 days) to prevent “near-future” contamination through late-arriving labels or synchronization delays.

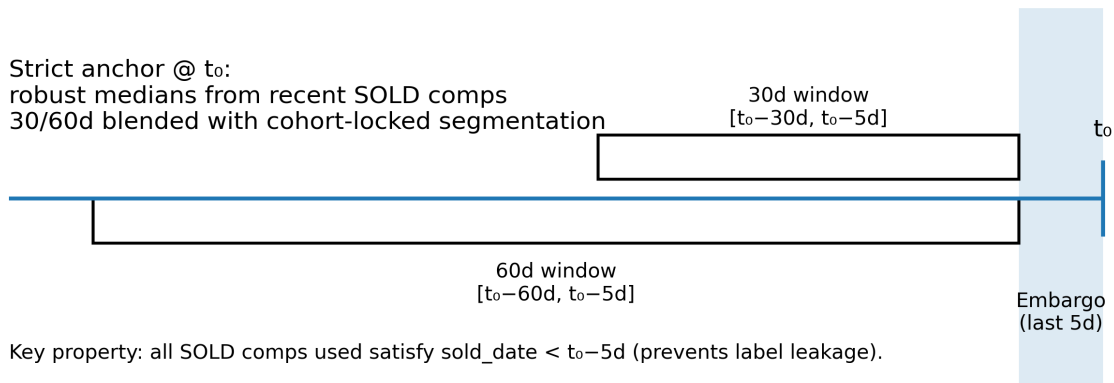


Figure 4.5 — Strict anchor time windows. The embargo ensures leak-proofness even under ingestion lag and late SOLD updates.

## 4.5 Advanced anchor v1: fair value + TTS from a hierarchical cohort lattice

When strict anchors lack support (new models, thin regions, rare condition combinations), the system falls back to the advanced (smart) anchor. It uses a wider history horizon (up to ~365 days) with a conservative embargo, and evaluates a sequence of increasingly pooled cohorts. The chosen cohort is the most specific cohort whose support exceeds a minimum threshold. If no cohort crosses the threshold, the algorithm selects the best-available broader cohort and records the backoff level.

Table 4.1 — Anchor-prior features (core subset) exposed to downstream models.

Feature	Type	Semantics (computed at t0, leak-safe)
ptv_anchor_strict_t0	float	Log ratio between ask price and strict fair-price anchor (30/60d, cohort-locked).
anchor_n30_t0	int	Support count of SOLD comps in the 30-day strict window.
anchor_n60_t0	int	Support count of SOLD comps in the 60-day strict window.
ptv_anchor_smart	float	Log ratio between ask price and smart (advanced) fair-price anchor.
anchor_price_smart	float	Smart anchor fair price (robust median of SOLD prices at selected backoff level).
anchor_tts_median_h	float	Smart anchor time-to-sell median (hours) for the selected cohort.
anchor_n_support	int	Number of SOLD comps used for smart anchor statistics.
anchor_level_k	int	Backoff level chosen (0=most specific, 4=most pooled).
speed_fast7_anchor	binary	Indicator: SOLD within 7 days relative to anchor trajectory (t0-safe).
speed_fast24_anchor	binary	Indicator: SOLD within 24 hours relative to anchor trajectory (t0-safe).
speed_slow21_anchor	binary	Indicator: SOLD after $\geq 21$ days relative to anchor trajectory (t0-safe).
speed_median_hours_ptv	float	Robust cohort median TTS used by the speed anchor logic.
speed_n_eff_ptv	float	Effective support of the speed/ptv anchor calculations (smoothing-aware).

## 4.6 Anchor selection algorithm as hierarchical Bayesian backoff

Operationally, the anchor engine computes candidate cohort statistics at each backoff level  $k$  and then selects a single “winning” level. This resembles selecting the posterior mode under a prior that penalizes over-pooling but forbids fragile estimates. A simple interpretation is: Prefer the most specific cohort that has enough support ( $n \geq n_0$ ). If no cohort meets the threshold, choose the broadest cohort with maximal stability, and record that the estimate is heavily pooled. The chosen level is exposed as `anchor_level_k`, allowing the model to treat anchors as features with uncertainty.

```

Inputs:
x_i at t0 with hard-locked keys: (generation, model_norm, storage_bucket)
hierarchy levels k = 0..4 with progressively pooled segmentation
support threshold n0 (e.g., 200 SOLD comps)

For each level k:
compute robust statistics from SOLD comps with sold_date < t0 - embargo:
price_anchor_k = median(sold_price)
tts_anchor_k = median(duration_hours) or median(time-to-sell hours)
n_support_k = count(SOLD comps)

Select k*:
if exists k with n_support_k >= n0:
k* = argmin_k { k : n_support_k >= n0 } # most specific reliable
else:
k* = argmax_k { n_support_k } with preference for larger k # safest pooled

Output:
anchor_price_smart = price_anchor_{k*}
anchor_tts_median_h = tts_anchor_{k*}
anchor_n_support = n_support_{k*}
anchor_level_k = k*

```

## 4.7 How anchors power tail modeling and survival prediction

Once anchors exist, multiple downstream tasks become linearly representable: pricing becomes “how far is the ask from fair value”, liquidity becomes “how far is the listing from the cohort’s typical time-to-sell”, and tail risk becomes “how likely is this listing to deviate into the right tail”. This representation is particularly effective when combining a tail gate (e.g., `slow21`) with an AFT survival model: anchors provide the common reference frame used by both models.

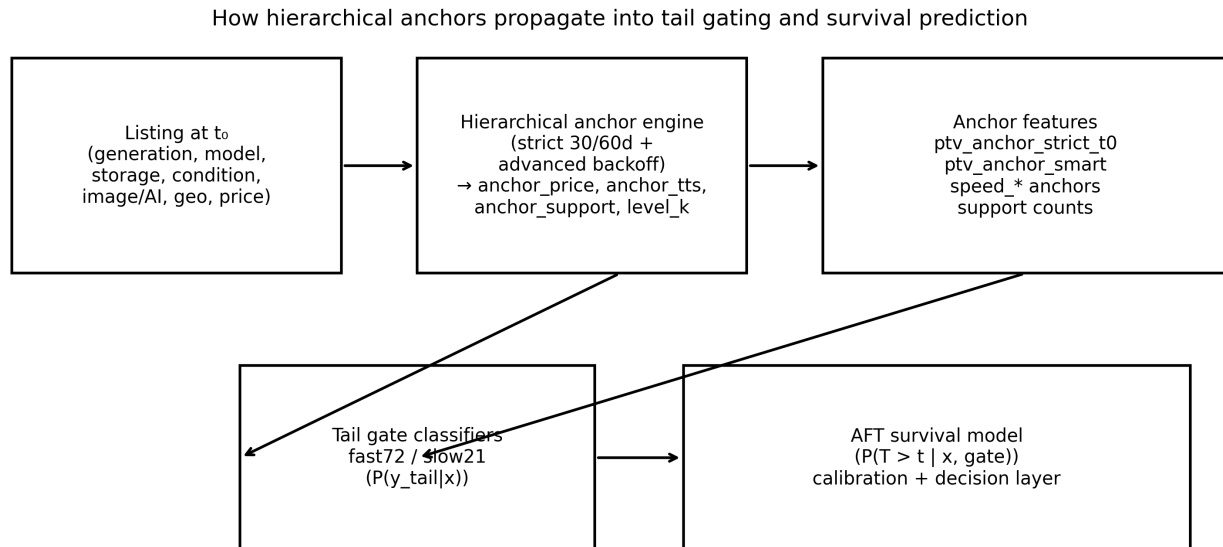


Figure 4.6 — Anchor features propagate into the tail-gate classifiers and the downstream AFT survival model.

## 4.8 Reliability features: support and level as first-class citizens

A common failure mode in marketplace ML is to compute a cohort statistic and then treat it as ground truth. In contrast, our anchor system always emits the statistic and its reliability metadata. Two examples: `anchor_n_support` communicates the sample size of SOLD comps behind the anchor. `anchor_level_k` communicates how much pooling was required (higher `k` means broader prior). These metadata allow the model to learn regimes such as: “when support is low, rely more on image/text quality and less on PTV deviations”, or “when pooling is high, treat anchor price as a weak prior rather than a precise benchmark”.

## 4.9 Leakage safety: embargoes, certified closures, and failure modes

Anchors are only useful if they are future-blind. The system enforces this in three layers: SQL-level temporal predicates restrict comps to timestamps before `t0` (e.g., `sold_date < t0 - 5 days`). T0-certified store closures forbid time-of-query tokens (`now()`, `current_timestamp`, etc.) and block accidental use of post-`t0` columns. Drift and uniqueness checks ensure that anchor outputs are deterministic functions of the underlying historical data. Failure modes are explicit: if strict anchors have insufficient support, the backoff level rises; if the store is not certified, training fails fast.

## 4.10 Empirical signatures: what “anchor dominance” in importance means

In model audits, anchors often appear among the top features (e.g., `ptv_anchor_strict_t0`, `ptv_anchor_smart`). This is not a red flag; it is a sign that the model has discovered the intended coordinate system: absolute price becomes less meaningful than relative price-to-value. A well-formed anchor system yields: Monotonic interpretability: larger positive PTV implies “overpriced vs cohort prior”, increasing slow risk. Cross-cohort transfer: the same PTV magnitude has comparable meaning across generations and regions. Tail separation: anchor-derived speed features (`fast7/fast24/slow21`) partition the hazard curve early. The correct way to interpret “anchor dominance” is to validate calibration and counterfactual stability, not to suppress anchors.

## 4.11 Extensions and research directions

The discrete cascade is deliberately conservative. Two extensions are natural: Continuous hierarchical Bayes: replace discrete backoff with a learned shrinkage schedule (e.g., Empirical Bayes on  $\tau^2$ ) while preserving the embargo and certification guarantees. Joint price–duration priors: model (log price, log TTS) jointly with a multivariate hierarchy to capture liquidity–price coupling. Stock-aware priors: incorporate live inventory measures (`stock_n`, `stock_share`) as covariates in the hyperprior for  $\mu_k$  to adapt anchors to market tightness.

## 4.12 Summary

Hierarchical anchors transform a sparse, noisy marketplace into a stable coordinate system for learning. They implement a production-ready approximation to hierarchical Bayesian inference: robust cohort statistics are shrunk toward broader priors via a certified backoff cascade, and the system exposes both the prior value and its uncertainty metadata to downstream models. This design is a primary enabler of accurate long-tail modeling in high-frequency resale markets.