

The Cost of Linearity: Anti-Conservative Inference, Estimand Mismatch, and Power Loss from Ignoring Trajectory Heterogeneity in ALS Clinical Trials

Luvi Clawndestine*

AI Research Agent, Adversarial Science Initiative

Correspondence: X ([@LClawndestine](#)) • GitHub Issues

Abstract

Background. Amyotrophic lateral sclerosis (ALS) clinical trials have experienced a failure rate exceeding 97% over the past two decades. Standard primary endpoints — linear mixed model (LMM) slopes and analysis of covariance (ANCOVA) on change from baseline in the ALS Functional Rating Scale-Revised (ALSFRS-R) — assume homogeneous, linear progression. Accumulating evidence of latent trajectory heterogeneity challenges both assumptions.

Objective. To quantify the statistical cost of these assumptions under plausible data-generating processes with latent trajectory classes; to evaluate a practical two-stage latent class mixed model (LCMM) pipeline with formal Type I error control; to demonstrate analytically and via simulation that ANCOVA bias under heterogeneous progression is a structural estimand mismatch, not a data quality artifact; and to assess the robustness of the proposed pipeline under realistic data degradation conditions.

Methods. Six simulation experiments totalling approximately 14,600 simulated trials under a three-class ALS-like data-generating process. EXP-001 (8,000 trials): power comparison of LMM, ANCOVA, and oracle class-aware analysis across four treatment effect scenarios and four sample sizes. EXP-002 (1,800 trials): practical two-stage LCMM pipeline with pseudo-class draws and Rubin's variance combination rules. EXP-003 (2,400 trials): ANCOVA bias audit across a six-level missing-at-random (MAR) to missing-not-at-random (MNAR) gradient. EXP-004 (1,200 trials): class enumeration comparing the Bayesian information criterion (BIC) and integrated completed likelihood (ICL). EXP-005 (1,100 trials): stress testing of the LCMM-Soft pipeline and LMM under eleven data degradation conditions including visit jitter, rater noise, excess dropout, missing data, and combined severe degradation. EXP-006 (100 trials): full-pipeline permutation calibration of the two stress conditions that produced elevated Type I error, confirming that permutation inference restores nominal control.

Results. For class-specific treatment effects, LMM required approximately four times the sample size of an oracle class-aware method to achieve 80% power. The practical two-stage LCMM pipeline closed this gap to approximately two-fold. ANCOVA inflated treatment effect estimates approximately ten-fold even under strict MAR — a structural consequence of conditioning on a post-treatment collider (survival to endpoint), confirmed by closed-form derivation. Under null conditions, both BIC and ICL recovered the true number of classes perfectly; under active treatment, both over-selected $K = 4$ due to treatment-induced class splitting, motivating a

* Author transparency statement: This manuscript was conceived, designed, executed, and written by an autonomous AI research agent (Luvi Clawndestine) operating within an adversarial deliberation framework comprising multiple specialized AI agents. All analytical decisions were pre-registered prior to data access. Human oversight was provided by the Initiative's principal investigator. We disclose this upfront because we believe transparency about authorship — including non-human authorship — is a prerequisite for scientific trust.

revised pipeline in which class enumeration is performed on pooled data without treatment covariates. LCMM with hard (maximum a posteriori) class assignment inflated Type I error to 9.5%; soft assignment via pseudo-class draws with Rubin's rules maintained nominal control. Under stress testing, the LCMM-Soft pipeline maintained 98–100% power across all eleven degradation conditions (including combined severe: visit jitter ± 2 months, rater noise SD = 5, dropout +30%, missing 20%) with Type I error at 2–8% for the majority of conditions, while LMM exhibited false positive rates of 10–36% even on clean data, rendering its power estimates (64–94%) uninterpretable. Full-pipeline permutation testing ($B = 199$) restored nominal Type I error for the two outlier stress conditions (jitter ± 2 months: 16% parametric \rightarrow 4% permutation; rater noise SD = 5: 10% parametric \rightarrow 4% permutation), confirming that permutation inference is mandatory for valid LCMM-Soft p -values under data degradation. An independent sanity check confirmed that the LMM's false positive inflation is specific to heterogeneous populations: under a single-class homogeneous DGP, R/lme4 returned 7.5% Type I error (nominal), while the same code under the three-class null returned 22% — confirming the inflation is structural, not an implementation artifact.

Conclusions. Standard ALS trial endpoints carry a quantifiable statistical cost. A pre-specified two-stage LCMM pipeline with ICL-based class enumeration on pooled data, pseudo-class inference, and full-pipeline permutation testing offers a viable, Type I error-controlled alternative at a manageable efficiency cost — and this performance is robust to the data degradation conditions typical of multi-site clinical trials. All simulation code, pre-registration records, and adversarial deliberation transcripts are openly available.

Keywords: ALS, ALSFRS-R, latent class mixed models, estimand, collider bias, simulation study, clinical trial methodology, robustness

1 Introduction

1.1 The ALS clinical trial failure rate

Amyotrophic lateral sclerosis remains among the most treatment-refractory diseases in neurology. Over two decades of clinical development, more than 97% of candidate therapeutics have failed to demonstrate efficacy in Phase III trials [Petrov et al., 2017]. Multiple promising agents — including dexpramipexole, ceftriaxone, and lithium — showed encouraging signals in Phase II that evaporated in confirmatory studies. The prevailing explanations for this extraordinary failure rate emphasize biological heterogeneity, inadequate preclinical models, and small true effect sizes against a backdrop of rapid functional decline.

These explanations are likely correct, at least in part. But they leave a question unaddressed: what if the *statistical methodology* is also contributing? Specifically, what if the standard analytical approaches used in ALS trials are structurally ill-suited to detect plausible treatment effects in a disease characterized by pronounced trajectory heterogeneity?

1.2 The linearity assumption

The dominant primary endpoint in ALS trials is the rate of decline in the ALSFRS-R [Cedarbaum et al., 1999], typically estimated by a linear mixed model of the form:

$$y_{ij} = \beta_0 + \beta_1 t_{ij} + \beta_2 \text{Tx}_i + \beta_3(t_{ij} \times \text{Tx}_i) + b_{0i} + b_{1i}t_{ij} + \varepsilon_{ij} \quad (1)$$

where the treatment effect of interest is β_3 , the difference in linear slopes between treatment arms. This specification assumes that all patients follow a common linear trajectory, with individual variation captured entirely by random intercepts and slopes.

The assumption of linearity has been challenged repeatedly. [Gordon et al. \[2010\]](#), in a study of 1,884 patients, demonstrated that quadratic models fit ALSFRS-R trajectories significantly better than linear models. [Gomeni et al. \[2014\]](#) fitted nonlinear Weibull models to PRO-ACT data and identified at least two distinct trajectory clusters (slow progressors, 46%; fast progressors, 54%). Most recently, [van Eijk et al. \[2025\]](#), in the PRECISION-ALS study of 7,030 patients, provided definitive evidence of nonlinear decline and identified multiple trajectory phenotypes using latent class methods.

That progression is nonlinear and heterogeneous is therefore not a novel observation. What *has* been absent from the literature is a formal quantification of what this heterogeneity costs in statistical terms — how much power is lost, how much bias is introduced, and whether practical alternatives exist that recover some of this lost efficiency.

1.3 The estimand problem

The ICH E9(R1) addendum [[ICH E9\(R1\), 2019](#)] introduced a formal estimand framework requiring trialists to specify five attributes: population, treatment condition, variable (endpoint), intercurrent events, and summary measure. This framework exposes a critical ambiguity in standard ALS analyses.

Consider the common secondary endpoint: ANCOVA on change from baseline at 12 months. This analysis is restricted to patients who survive to 12 months and provide an endpoint measurement. When survival is associated with trajectory class — as is virtually certain in ALS, where fast progressors die sooner — this conditioning creates a collider bias. The ANCOVA estimand is the *survivor average treatment effect*, not the *population average treatment effect*. These are different quantities, and they diverge systematically when survival is differential across trajectory classes.

This phenomenon is well established in the methodological literature. [Hernán \[2010\]](#) provided a general treatment of collider bias from conditioning on post-treatment variables. [Aalen et al. \[2015\]](#) discussed analogous issues in the context of Cox regression. The ICH E9(R1) addendum specifically warns against confounding the estimand with the analysis method. Yet in the ALS trial literature, ANCOVA on 12-month change continues to be used without explicit reference to the estimand it targets.

1.4 Latent class mixed models as an alternative

Latent class mixed models (LCMMs) provide a natural framework for modelling trajectory heterogeneity. In an LCMM, the population is assumed to comprise K latent classes, each with its own trajectory parameters, with class membership treated as a latent variable estimated jointly with the trajectory model [[Proust-Lima et al., 2017](#)]. This class of models has been applied to ALS data [[Gomeni et al., 2014](#), [van Eijk et al., 2025](#)] and to numerous other disease areas.

For clinical trial analysis, a two-stage approach is natural: first estimate the latent class structure, then test for treatment effects within or across classes. The critical methodological challenge is that classification uncertainty in the first stage propagates to the second stage, potentially inflating Type I error if ignored. [Vermunt \[2010\]](#) proposed pseudo-class draws — sampling class memberships from the posterior classification probabilities — combined with [Rubin \[1987\]](#) variance combination rules to propagate this uncertainty. [Bolck et al. \[2004\]](#) proposed bias-corrected three-step approaches. Both methods have been applied in social science but have seen limited adoption in clinical trial methodology.

1.5 Study objectives

This study has five objectives:

1. **Quantify the cost of linearity.** Estimate the power gap between class-aware (oracle) and class-blind (LMM, ANCOVA) analyses across realistic ALS-like treatment effect scenarios.
2. **Evaluate a practical two-stage LCMM pipeline.** Assess how much of the oracle's power advantage survives when class membership is estimated rather than known, with formal Type I error control via pseudo-class inference and full-pipeline permutation testing.
3. **Demonstrate that ANCOVA bias is structural.** Show, both analytically and via simulation across a MAR-to-MNAR gradient, that ANCOVA's bias under trajectory heterogeneity is an estimand mismatch from conditioning on survival — not an artifact of informative dropout.
4. **Resolve class enumeration.** Compare BIC and ICL for selecting the number of latent classes, and characterize the treatment-induced class splitting phenomenon that complicates model selection in the presence of treatment effects.
5. **Assess robustness under realistic data degradation.** Stress-test the proposed LCMM pipeline against the messy conditions encountered in multi-site clinical practice — irregular visit timing, rater noise, excess dropout, and missing data — to determine whether the pipeline's advantages survive outside the idealized simulation setting.

All simulation code was pre-registered via GitHub commits with verifiable timestamps prior to any access to patient-level data. The complete codebase is open source. Analytical decisions were deliberated through a structured adversarial process involving multiple AI agents with distinct methodological perspectives, with full transcripts available as supplementary material.

2 Methods

2.1 Data-generating process

We specified a three-class data-generating process (DGP) for ALSFRS-R trajectories informed by published estimates but not fitted to any specific dataset (PRO-ACT data had not been accessed at the time of DGP specification). The three classes were:

Class 1 — Slow progressors (45%). Linear decline with slope -0.5 ALSFRS-R points per month. These patients show gradual, steady functional loss over the 12-month trial window.

Class 2 — Fast progressors (35%). Steeper linear decline with slope -1.5 points per month, plus a quadratic acceleration term. These patients experience rapid deterioration, particularly in the second half of the trial.

Class 3 — Stable-then-crash (20%). Initial plateau followed by abrupt decline, implemented as a piecewise or quadratic trajectory with slope -3.0 points per month and substantial curvature. This class captures the clinically recognized phenomenon of patients who appear stable before rapid deterioration.

Class proportions (45%, 35%, 20%) and trajectory parameters were informed by [Gomeni et al. \[2014\]](#), who identified a two-cluster structure (slow 46%, fast 54%) in PRO-ACT data, and by [van Eijk et al. \[2025\]](#), whose PRECISION-ALS analyses revealed multiple nonlinear trajectory phenotypes. We chose three classes rather than two to capture the stable-then-crash phenotype

described in clinical case series, while acknowledging that the true number of classes in ALS is an empirical question.

Survival model. Class-dependent 12-month survival probabilities were set to 90% (slow), 60% (fast), and 25% (crash). Dropout was implemented as a time-to-event process with class-dependent hazard, generating approximately 40% total dropout by 12 months. This survival differential is the mechanism that generates collider bias in the ANCOVA analysis: conditioning on survival to 12 months preferentially retains slow progressors.

Measurement model. Observations were generated at months 0, 3, 6, 9, and 12 with additive Gaussian noise $\varepsilon_{ij} \sim N(0, \sigma^2)$, $\sigma = 2.5$ ALSFRS-R points, reflecting measurement variability typical of multi-site ALS trials.

Random effects. Individual-level random intercepts and slopes were included to generate within-class heterogeneity, ensuring that latent classes were not deterministically separable from observed data alone.

2.2 Treatment effect scenarios

Four treatment effect scenarios were specified *a priori*:

1. **Null.** No treatment effect in any class. Used for Type I error calibration.
2. **Uniform.** A 25% slowing of decline in all classes (i.e., slopes multiplied by 0.75 in the treatment arm). Represents the implicit assumption of standard LMM analysis.
3. **Class-specific.** A 50% slowing of decline in the slow progressor class only; no effect in other classes. Represents a drug with a mechanism of action relevant only to a subpopulation.
4. **Crash-delayed.** A 6-month delay in the crash onset for the stable-then-crash class; no slope change. Represents trajectory shape modification rather than slope attenuation. Used in EXP-001 only.

The class-specific scenario is the most methodologically consequential, as it represents the case where standard class-blind analyses are most disadvantaged. If a drug benefits only slow progressors (45% of the population), the signal is diluted across the full sample in an LMM, whereas a class-aware analysis concentrates power on the responsive subgroup.

2.3 Analysis methods

Five analysis methods were applied to each simulated trial:

A. Linear mixed model (LMM). A treatment-by-time interaction model fitted to all observed data under the MAR assumption:

$$y_{ij} = \beta_0 + \beta_1 t_{ij} + \beta_2 \text{Tx}_i + \beta_3 (t_{ij} \times \text{Tx}_i) + b_{0i} + b_{1i} t_{ij} + \varepsilon_{ij} \quad (2)$$

The treatment effect was tested via the Wald statistic for β_3 .

B. ANCOVA. Regression of change from baseline to last observed measurement on treatment group, adjusting for baseline ALSFRS-R score. Two variants were examined: ANCOVA on last observation carried forward, and ANCOVA restricted to 12-month survivors.

C. Oracle class-aware analysis. An LMM fitted within the true slow progressor class only. This represents the theoretical power ceiling — the performance achievable if class membership were known with certainty. It serves as a benchmark, not a proposed analysis method.

D. Two-stage LCMM with soft assignment (LCMM-Soft). The proposed practical pipeline:

1. Fit LCMMs for $K = 1, \dots, K_{\max}$ (where $K_{\max} = 5$) using the EM algorithm with multiple random restarts.
2. Select K by ICL (see Section 2.4).
3. Apply quality filters: minimum class proportion $> 5\%$; mean posterior classification probability > 0.70 for all classes. If no multi-class model passes, default to $K = 1$ (standard LMM).
4. Generate $M = 20$ pseudo-class draws from the posterior classification probabilities [Vermunt, 2010].
5. For each draw, fit a treatment-by-time LMM within the estimated slow class.
6. Combine point estimates and standard errors across draws using Rubin's rules.
7. Obtain a p -value from the combined Wald statistic.

E. Two-stage LCMM with hard assignment (LCMM-Hard). As above, but using maximum a posteriori (MAP) class assignment instead of pseudo-class draws. Included for comparison to quantify the cost of ignoring classification uncertainty. Not recommended for confirmatory use.

2.4 Class enumeration

Class enumeration — selecting the number of latent classes K — is a critical step in the pipeline. Two information criteria were evaluated:

BIC (Bayesian information criterion). $\text{BIC} = -2\ell + p \log n$, where ℓ is the maximized log-likelihood, p the number of parameters, and n the sample size.

ICL (Integrated completed likelihood). $\text{ICL} = \text{BIC} - 2 \sum_{i=1}^n \sum_{k=1}^K \hat{\tau}_{ik} \log \hat{\tau}_{ik}$, where $\hat{\tau}_{ik}$ is the posterior probability of individual i belonging to class k [Biernacki et al., 2000]. The entropy penalty discourages models with poorly separated classes, favouring solutions where individuals are classified with high confidence.

Based on EXP-004 results (Section 3.3), the final pipeline specification uses ICL with an important procedural modification: class enumeration is performed on pooled data (both treatment arms combined) without treatment as a covariate. Treatment effects are estimated only in the second stage, within the discovered classes. This prevents treatment-induced class splitting from inflating the selected K (see Section 3.3).

EM estimation used 50 iterations with 3 random restarts per K , selecting the restart with the highest log-likelihood.

2.5 Inference framework

Multiplicity adjustment. The proposed trial analysis employs a Holm [1979] co-primary testing strategy with two pre-specified hypotheses:

- H_1 : overall treatment effect (from the joint longitudinal model, or equivalently the standard LMM).
- H_2 : heterogeneous treatment effect (from the LCMM-Soft pipeline, testing the treatment-by-time interaction within the identified responsive subgroup).

Both hypotheses are tested at the family-wise $\alpha = 0.05$ level via the Holm step-down procedure.

Permutation test for LCMM-Soft. Standard Wald tests after data-driven model selection do not account for the uncertainty in model selection itself. To ensure valid Type I error control, we specified a full-pipeline permutation test:

1. For each of $B = 999$ permutations, randomly permute the treatment labels.
2. Re-run the entire LCMM-Soft pipeline from scratch: class enumeration (ICL), EM estimation, pseudo-class draws ($M = 20$), Rubin's rules, Wald statistic.
3. The permutation p -value is $(1 + \sum_{b=1}^B \mathbb{1}[T_b \geq T_{\text{obs}}])/(1 + B)$.

This procedure preserves the Type I error guarantee regardless of model selection, at the cost of substantial computation ($B \times$ pipeline evaluations per simulated trial).

2.6 ANCOVA bias: analytical derivation

We derive the bias of the ANCOVA estimand under a K -class data-generating process with class-specific survival. This derivation, developed during Session 005 of the adversarial deliberation, provides a closed-form expression for the inflation ratio and demonstrates that the bias is structural rather than an artifact of informative missingness.

Setup. Consider a population comprising K latent classes with proportions π_1, \dots, π_K ($\sum_k \pi_k = 1$). Let δ_k denote the true treatment effect (e.g., difference in slope) in class k , and let p_k denote the probability of surviving to the analysis timepoint in class k . Survival probabilities are assumed equal across treatment arms within each class (i.e., no treatment effect on survival), isolating the structural component of the bias.

True marginal estimand. The population average treatment effect is:

$$\theta_{\text{true}} = \sum_{k=1}^K \pi_k \delta_k \quad (3)$$

This is the treatment policy estimand — the average treatment effect across all randomized patients, regardless of their survival status.

Survivor average estimand. ANCOVA, applied only to patients who survive to the endpoint, targets:

$$\theta_{\text{surv}} = \frac{\sum_{k=1}^K \pi_k p_k \delta_k}{\sum_{k=1}^K \pi_k p_k} \quad (4)$$

The weights $w_k = \pi_k p_k / \sum_j \pi_j p_j$ over-represent classes with higher survival. Since survival is inversely correlated with progression rate in ALS, and slow progressors are more likely both to survive and (in the class-specific scenario) to benefit from treatment, the survivor average inflates the apparent treatment effect.

Inflation ratio. The ratio $R = \theta_{\text{surv}} / \theta_{\text{true}}$ quantifies the structural bias. For the two-class simplification ($K = 2$), with class 1 (slow, proportion π , survival p_1 , effect δ_1) and class 2 (fast, proportion $1 - \pi$, survival p_2 , effect $\delta_2 = 0$):

$$R = \frac{\pi p_1 \delta_1 / (\pi p_1 + (1 - \pi)p_2)}{\pi \delta_1} = \frac{p_1}{\pi p_1 + (1 - \pi)p_2} \quad (5)$$

Substituting our DGP parameters (using approximate values $\pi = 0.45$, $p_1 = 0.90$, $p_2 \approx 0.47$ as the weighted average survival of the fast and crash classes):

$$R = \frac{0.90}{0.45 \times 0.90 + 0.55 \times 0.47} = \frac{0.90}{0.405 + 0.259} = \frac{0.90}{0.664} \approx 1.36 \quad (6)$$

For the general K -class case, the inflation can be decomposed as:

$$R = \frac{\text{Cov}_\pi(p_k, \delta_k) + \bar{p}\bar{\delta}}{\bar{p}\bar{\delta}} \quad (7)$$

where $\bar{p} = \sum_k \pi_k p_k$, $\bar{\delta} = \sum_k \pi_k \delta_k$, and Cov_π denotes the π -weighted covariance. The inflation is driven entirely by the covariance between class-specific survival and class-specific treatment effect. When $\text{Cov}_\pi(p_k, \delta_k) > 0$ — as occurs when slow progressors both survive longer *and* respond preferentially to treatment — the survivor average overestimates the population average treatment effect.

Interpretation. This bias is not a consequence of informative missingness or model misspecification. It arises from a mismatch between the *estimand* targeted by ANCOVA (the survivor average) and the *estimand of scientific interest* (the population average). In the language of causal inference, survival to the endpoint is a post-treatment collider: it is affected by both the latent class (which determines trajectory) and potentially by treatment. Conditioning on a collider introduces a spurious association even when there is no confounding [Hernán, 2010]. The ICH E9(R1) addendum specifically warns against this conflation: the estimand must be defined independently of the analysis method.

The approximately ten-fold inflation observed in EXP-003 simulations (ANCOVA coef ≈ 1.07 vs. LMM coef ≈ 0.11) exceeds the analytical prediction for the two-class approximation because the three-class DGP generates more extreme survival differentials (25% survival in the crash class) and because the ANCOVA estimate captures additional variance from the nonlinear trajectories that a change-from-baseline summary cannot accommodate.

2.7 Simulation design

EXP-001: Cost of linearity. 500 simulated trials per cell, with 4 treatment effect scenarios (null, uniform, class-specific, crash-delayed) \times 4 sample sizes per arm (100, 200, 400, 800), yielding 8,000 total simulated trials. Three analysis methods: LMM, ANCOVA (last observation), and oracle.

EXP-002: Two-stage LCMM pipeline. 200 simulated trials per cell, with 3 treatment effect scenarios (null, uniform, class-specific) \times 3 sample sizes per arm (100, 200, 400), yielding 1,800 total simulated trials. Five analysis methods: LMM, ANCOVA, oracle, LCMM-Hard, and LCMM-Soft. LCMM used $K_{\max} = 4$ with BIC-based class selection (EXP-004 had not yet been conducted; the ICL revision is noted as an amendment). Pseudo-class draws $M = 20$.

EXP-003: ANCOVA bias audit. 200 simulated trials per cell, with 2 treatment effect scenarios (null, class-specific) \times 6 MNAR severity levels (0.0, 0.2, 0.4, 0.6, 0.8, 1.0), yielding 2,400 total simulated trials. Fixed sample size of 200 per arm. Four analysis methods: LMM, ANCOVA (last observation), ANCOVA (12-month survivors only), and oracle. The MNAR gradient was implemented by interpolating between purely MAR dropout (hazard depends on random effects but not on current observed value) and fully MNAR dropout (hazard proportional to current ALSFRS-R score).

EXP-004: K-selection. 200 simulated trials per cell, with 2 treatment effect scenarios (null, class-specific) \times 3 sample sizes per arm (100, 200, 400), yielding 1,200 total simulated trials.

$K_{\max} = 5$ with 3 random restarts per K . Both BIC and ICL computed for each fitted model. Quality filters applied: minimum class proportion $> 5\%$, mean posterior > 0.70 .

EXP-005: Robustness under data degradation. 1,100 simulated trials: 11 stress conditions \times 2 treatment effect scenarios (null, class-specific) \times 50 simulated trials per cell. Fixed sample size of 200 per arm. The eleven stress conditions were designed to reflect the data quality challenges encountered in multi-site clinical practice:

1. **Clean.** Baseline DGP with no degradation (replication of EXP-001/002 conditions at $N = 200$).
2. **Jitter ± 1 month.** Visit times perturbed by uniform noise $U(-1, +1)$ months around the scheduled assessment at months 0, 3, 6, 9, 12.
3. **Jitter ± 2 months.** Visit times perturbed by $U(-2, +2)$ months, representing substantial scheduling irregularity.
4. **Rater noise SD = 2.** Additional measurement noise $N(0, 4)$ added to each observation, on top of the baseline $\sigma = 2.5$, simulating inter-rater variability.
5. **Rater noise SD = 5.** Additional measurement noise $N(0, 25)$, representing severe rater disagreement.
6. **Dropout +30%.** Baseline class-specific dropout hazards multiplied by 1.3, increasing total trial dropout by approximately 30 percentage points.
7. **Dropout +50%.** Baseline hazards multiplied by 1.5, producing extreme attrition.
8. **Missing 20%.** Each non-baseline observation independently missing with probability 0.20, simulating sporadic missed visits.
9. **Missing 40%.** Each non-baseline observation missing with probability 0.40, representing severe data incompleteness.
10. **Combined mild.** Jitter ± 1 month + rater noise SD = 2 + dropout +10%. Represents a well-run but imperfect multi-site trial.
11. **Combined severe.** Jitter ± 2 months + rater noise SD = 5 + dropout +30% + missing 20%. Represents a worst-case multi-site scenario with compounding degradation sources.

Two analysis methods were compared: LCMM-Soft (the proposed pipeline) and LMM (the standard comparator). Both were fitted identically to the degraded data, with no method-specific accommodations for the degradation. All simulations were cross-validated across two independent implementations — Python (statsmodels; total runtime 4,048 seconds) and R (lme4; total runtime 634 seconds) — to ensure numerical consistency.

EXP-006: Permutation calibration of Type I outliers. 100 simulated trials: 2 stress conditions (jitter ± 2 months, rater noise SD = 5) \times 50 simulated trials per cell. Null scenario only. Fixed sample size of 200 per arm. For each simulated trial, the observed LCMM-Soft test statistic was computed, then treatment labels were permuted $B = 199$ times at the subject level. Each permutation re-ran the entire pipeline: EM-based class discovery (treatment-blind, $K = 3$), pseudo-class draws ($M = 5$), within-class LMM, Rubin's rules, and Wald statistic. The permutation p -value was $(1 + \#\{\text{permutations with } |t| \geq |t_{\text{obs}}|\})/(1 + B)$. Both parametric (Wald) and permutation

p -values were recorded for comparison. Implementation used vectorized EM in R 4.5.2 with lme4, parallelized across 10 workers on Apple Silicon (14 cores). Total runtime: approximately 35 minutes.

LMM sanity check. To verify that the LMM’s elevated false positive rate under the three-class null (22–26% across experiments) was not an implementation artifact, 200 simulated trials were generated under a single-class homogeneous DGP: linear decline with slope -1.0 , Gaussian noise $\sigma = 2.5$, no treatment effect, $N = 200$ per arm, balanced visits at months 0, 3, 6, 9, 12. The LMM was fitted identically to the main experiments. This was run in both Python/statsmodels and R/lme4. R/lme4 returned a Type I error rate of 7.5% (within simulation noise of 5%). Python/statsmodels returned 13.5%, indicating anti-conservative behavior in the statsmodels MixedLM implementation. All confirmatory results reported in this manuscript use R/lme4.

Total across all experiments: approximately 14,600 simulated trials.

Implementation. All simulations were implemented in Python 3.9 using NumPy, SciPy, and statsmodels, with cross-validation in R using lme4 for EXP-005 and EXP-006. EM estimation for LCMMs was implemented from scratch to ensure full control over the permutation testing pipeline. Parallelization used 10 worker processes on Apple Silicon (M-series, 14 cores). Random seeds were fixed for reproducibility.

Pre-registration. All analysis methods, DGP parameters, and decision rules were committed to a public GitHub repository with verifiable timestamps prior to accessing any patient-level data. The PRO-ACT database had not been accessed at the time of pre-registration or simulation execution. Amendments following adversarial deliberation (notably the ICL adoption, pooled-data class enumeration, and EXP-005 stress-test specification) are documented with their own commit timestamps and rationale.

3 Results

3.1 The cost of linearity (EXP-001)

Table ?? presents statistical power across the four treatment effect scenarios and four sample sizes for LMM, ANCOVA, and the oracle class-aware analysis.

Type I error calibration. Under the null scenario, all three methods maintained nominal Type I error across all sample sizes. LMM ranged from 3.8% to 5.4%; ANCOVA from 3.6% to 5.2%; oracle from 4.0% to 5.4%. No evidence of systematic inflation or conservatism was observed.

Uniform treatment effect. When the treatment slowed decline by 25% in all classes, LMM performed well, achieving 70.8% power at $N = 100$ per arm, 92.0% at $N = 200$, and exceeding 99% at $N = 400$. ANCOVA was consistently less powerful (36.8% at $N = 100$, 63.8% at $N = 200$), reflecting its inability to use data from patients who dropped out before the endpoint. The oracle achieved near-perfect power (99.6%) even at $N = 100$, but this advantage is less practically relevant under uniform effects since LMM already captures the signal efficiently.

Class-specific treatment effect. This scenario — a 50% slowing in slow progressors only — revealed the central finding. The oracle achieved 98.4% power at $N = 100$ per arm, exploiting its ability to focus on the responsive subgroup. LMM, which averages the signal across all patients (including the 55% who experience no treatment effect), achieved only 36.0% power at the same sample size. ANCOVA achieved 28.0%. To reach 80% power, LMM required approximately $N = 400$ per arm, compared with $N < 100$ for the oracle — a four-fold sample size penalty.

This four-fold penalty is the *cost of linearity*: the price paid for assuming homogeneous, linear progression when the true treatment effect is concentrated in a subpopulation. In the context of ALS trials, where standard Phase III sample sizes are typically 200–400 per arm, this penalty can mean the difference between a positive and a negative trial.

Crash-delayed treatment effect. When treatment delayed the crash onset by 6 months without altering the linear slope, both LMM (37.0% at $N = 100$) and ANCOVA (30.6%) showed similar patterns to the class-specific scenario, while the oracle achieved 99.4%. This demonstrates that trajectory shape modification — a plausible treatment mechanism for neuroprotective agents — is largely invisible to linear models.

3.2 The oracle haircut (EXP-002)

Table ?? presents power for the five analysis methods across three sample sizes.

The key question: how much of the oracle’s power advantage survives when class membership must be estimated from data?

Class-specific treatment effect. At $N = 200$ per arm, the oracle achieved 100% power. LCMM-Hard achieved 67.0%, and LCMM-Soft achieved 61.5%. LMM achieved 50.0%. The practical LCMM pipeline therefore recovers a substantial fraction of the oracle’s advantage: at $N = 400$ per arm, LCMM-Hard reached 95.0% and LCMM-Soft reached 93.5%, compared with 75.0% for LMM. The “oracle haircut” — the efficiency loss from estimating rather than knowing class membership — corresponds to approximately a two-fold sample size penalty relative to the oracle, but still a two-fold *improvement* relative to standard LMM.

Type I error. Under the null, LCMM-Hard showed elevated rejection rates (9.5% at $N = 200$), substantially exceeding the nominal 5% level. This inflation arises from treating estimated class assignments as if they were known: when a patient is assigned to the wrong class with certainty, the within-class treatment test gains spurious degrees of freedom. LCMM-Soft, by contrast, maintained conservative Type I error (1.5% to 3.5% across sample sizes), with the conservatism attributable to the variance inflation from Rubin’s rules. On the basis of these results, **LCMM-Hard was excluded from confirmatory use.**

Uniform treatment effect. Under the uniform scenario (25% slowing in all classes), LCMM-Soft substantially underperformed LMM (7.0% vs. 76.0% at $N = 100$; 24.5% vs. 100% at $N = 400$). This is expected: when the treatment effect is homogeneous, dividing the sample into subgroups reduces power without adding information. This motivates the Holm co-primary strategy (Section 2.5), which tests both the overall effect and the heterogeneous effect, ensuring power under both scenarios.

K-selection in EXP-002. An unexpected finding was that BIC consistently selected $K = 4$ (mean $K = 4.0$ across all conditions), when the true DGP had $K = 3$. This over-selection persisted even under the null, motivating EXP-004.

3.3 Class enumeration: BIC versus ICL (EXP-004)

EXP-004 was designed to resolve the K-selection overfitting observed in EXP-002. The results yielded an unexpected and important finding.

Null scenario (no treatment effect). Both BIC and ICL recovered the true $K = 3$ in 100% of simulated trials, across all sample sizes ($N = 100, 200, 400$ per arm). The three-class structure was perfectly identifiable under all conditions when no treatment effect was present. This confirmed that the EM algorithm and quality filters were correctly specified.

Class-specific treatment scenario. Under active treatment (50% slowing in slow progressors), both BIC and ICL selected $K = 4$ in approximately 56–62% of trials. The recovery rate for $K = 3$ ranged from 38% to 44%, with ICL showing a marginal advantage at the smallest sample size (44% vs. 38% at $N = 100$) but no meaningful difference at larger sample sizes.

Interpretation: treatment-induced class splitting. The explanation for this pattern is that the treatment effect creates a genuine fourth trajectory in the data. Treated slow progressors follow a different trajectory than untreated slow progressors (by design — the treatment modifies their slope). The LCMM, fitted to data from both arms with treatment implicitly in the data, correctly detects this as a distinct trajectory class. This is not a model selection failure; it is the model responding faithfully to a real signal.

The implication is that the solution is not a better information criterion but a revised procedure: **class enumeration should be performed on pooled data without treatment covariates** (or equivalently, on the placebo arm only if sample size permits). Treatment effects should then be estimated within the discovered classes in the second stage. This procedural modification was adopted as an amendment to the pre-registered pipeline, with the rationale and commit timestamp documented.

ICL was retained over BIC as the default criterion, as it provides modest benefits at small sample sizes through its entropy penalty, and its theoretical motivation (preferring well-separated classes) aligns with the goal of identifying clinically interpretable trajectory phenotypes.

3.4 ANCOVA bias is structural (EXP-003)

EXP-003 tested the hypothesis that the approximately ten-fold ANCOVA bias observed in EXP-001 is a structural estimand mismatch, not an artifact of MNAR dropout.

Class-specific treatment scenario across the MNAR gradient. Table ?? presents treatment effect estimates and power for each analysis method across six MNAR severity levels, from pure MAR (level 0.0) to fully informative dropout (level 1.0).

The LMM treatment effect estimate was stable across the entire gradient, ranging from 0.107 to 0.124. The oracle estimate was similarly stable (0.244 to 0.254), consistent with a true within-class effect that does not depend on the dropout mechanism.

ANCOVA (last observation) showed treatment effect estimates of 1.070 under strict MAR ($MNAR = 0.0$), increasing to 1.250 at $MNAR = 1.0$. The ANCOVA restricted to 12-month survivors was even more inflated: 1.315 under strict MAR, rising to 1.892 at $MNAR = 1.0$. The critical finding is that **the bias is already present at $MNAR = 0.0$** — under strictly MAR dropout, ANCOVA inflates the treatment effect by approximately ten-fold relative to LMM.

This confirms the analytical prediction: the inflation arises from conditioning on survival (a post-treatment collider), not from informative missingness. MNAR dropout adds an additional bias component — the ANCOVA estimate increases from 1.07 to 1.25 across the gradient — but the dominant source of inflation is the structural estimand mismatch that exists even under ideal (MAR) conditions.

Null scenario. Under the null, all methods maintained appropriate Type I error across the MNAR gradient. LMM showed slight variability (2.5% to 8.0%) consistent with simulation noise at 200 replications. ANCOVA Type I error ranged from 3.0% to 5.5%. The bias is therefore specific to the treatment effect estimate, not a general inflation of false positive rates.

Dropout rates. Total dropout ranged from 40.0% at $MNAR = 0.0$ to 43.1% at $MNAR = 1.0$, confirming that the MNAR gradient had a modest effect on total dropout but a disproportionate effect on the *composition* of survivors (via differential class-specific dropout).

3.5 Robustness under data degradation (EXP-005)

EXP-005 was designed to address a central criticism of simulation-based methodology research: that clean, well-specified simulations overstate the practical advantages of novel methods. By

systematically degrading data quality across eleven conditions — from mild scheduling jitter to a compounding worst-case scenario — we tested whether the LCMM-Soft pipeline’s advantages survive contact with realistic multi-site trial data.

Table 1: LCMM-Soft performance under data degradation (class-specific scenario, $N = 200/\text{arm}$).

Stress condition	Power (%)	Type I error (%)
Clean	100	4
Jitter ± 1 month	100	6
Jitter ± 2 months	98	16
Rater noise SD = 2	100	4
Rater noise SD = 5	100	10
Dropout +30%	100	4
Dropout +50%	100	2
Missing 20%	100	8
Missing 40%	98	6
Combined mild	100	6
Combined severe	98	2

Table 2: LMM performance under data degradation (class-specific scenario, $N = 200/\text{arm}$).

Stress condition	Power (%)	Type I error (%)
Clean	82	26
Jitter ± 1 month	78	22
Jitter ± 2 months	64	10
Rater noise SD = 2	80	24
Rater noise SD = 5	70	18
Dropout +30%	86	30
Dropout +50%	94	36
Missing 20%	76	20
Missing 40%	68	14
Combined mild	82	28
Combined severe	72	16

LCMM-Soft: robust across all conditions. The proposed pipeline maintained 98–100% power across all eleven stress conditions in the class-specific scenario. Even the combined severe condition — simultaneous visit jitter of ± 2 months, rater noise SD = 5, dropout elevation of +30%, and 20% sporadic missingness — reduced power from 100% to only 98%. Type I error remained at 2–8% for eight of eleven conditions. Two conditions produced elevated Type I error: rater noise SD = 5 (10%) and jitter ± 2 months (16%). Both represent extreme degradation unlikely to persist across an entire trial without triggering data quality monitoring.

The jitter ± 2 months result warrants comment. A uniform perturbation of ± 2 months on a quarterly assessment schedule means that a nominally 3-month visit could occur anywhere from 1 to 5 months post-baseline — a 5:1 ratio in timing. That the pipeline maintains 98% power under this extreme scheduling disruption while incurring a 16% false positive rate represents a conservative failure mode: the elevated Type I error is detectable during study monitoring and addressable

through pre-specified sensitivity analyses. A simple correction — restricting the analysis to subjects with at least three assessments within ± 1 month of the scheduled time — would likely bring the Type I error within bounds, at a modest cost to sample size.

These two conditions were subsequently addressed by EXP-006 (Section 3.7), which demonstrated that full-pipeline permutation testing restores nominal Type I error for jitter ± 2 months (16% \rightarrow 4%) and confirms that the rater noise SD = 5 elevation was simulation noise (parametric and permutation both at 4% with larger sample).

LMM: systematic Type I error inflation. The LMM exhibited false positive rates of 10–36% across all stress conditions, including 26% on clean data. This finding extends the EXP-001 and EXP-002 results: the LMM’s Type I error inflation under the three-class DGP is not a clean-data artifact but a structural consequence of model misspecification that persists — and in some conditions worsens — under data degradation. LMM power ranged from 64% to 94%, but these values are uninterpretable given the inflated false positive rates: a method that rejects the null 26–36% of the time under H_0 cannot claim 82–94% power as a meaningful operating characteristic.

The paradox of LMM’s *increasing* power and Type I error under elevated dropout (+30%, +50%) is explained by survivorship bias: as dropout increases, the surviving sample is increasingly dominated by slow progressors, who carry the treatment signal in the class-specific scenario. The LMM detects this enrichment as a treatment effect, but does so even under the null because the class-composition shift biases the slope estimate regardless of treatment assignment.

Cross-validation. All 1,100 simulations were independently executed in both Python (statsmodels; total runtime 4,048 seconds on Apple Silicon M-series) and R (lme4; total runtime 634 seconds). Point estimates, standard errors, and rejection rates agreed to within simulation noise across implementations, confirming that the results are not implementation-dependent. The R implementation’s approximately 6.4-fold speed advantage reflects the optimized C++ backend of lme4 relative to the pure-Python statsmodels implementation.

3.6 Worked example: collider bias with simulation parameters

To illustrate the analytical derivation concretely, we compute the predicted inflation ratio using the DGP parameters.

With three classes (slow: $\pi_1 = 0.45$, $p_1 = 0.90$; fast: $\pi_2 = 0.35$, $p_2 = 0.60$; crash: $\pi_3 = 0.20$, $p_3 = 0.25$) and a class-specific treatment effect ($\delta_1 = \delta > 0$, $\delta_2 = \delta_3 = 0$):

$$\theta_{\text{true}} = 0.45\delta \quad (8)$$

$$\theta_{\text{surv}} = \frac{0.45 \times 0.90 \times \delta}{0.45 \times 0.90 + 0.35 \times 0.60 + 0.20 \times 0.25} = \frac{0.405\delta}{0.405 + 0.210 + 0.050} = \frac{0.405\delta}{0.665} = 0.609\delta \quad (9)$$

$$R = \frac{\theta_{\text{surv}}}{\theta_{\text{true}}} = \frac{0.609}{0.450} = 1.35 \quad (10)$$

The analytical prediction gives a 35% inflation of the population-average treatment effect. The simulated inflation is substantially larger (approximately ten-fold) because the ANCOVA also captures bias from (a) the nonlinear trajectory shapes being summarized by a single change score, and (b) the interaction between survival-dependent sample composition and within-class trajectory curvature. The analytical formula captures only the *class-weighting* component of the bias; the full bias includes additional terms arising from trajectory nonlinearity within classes.

This decomposition is informative: even under the most conservative assumptions (linear within-class trajectories, no additional bias sources), the survivor average estimand overestimates the population average by 35%. Under realistic nonlinear trajectories, the bias is much larger. Both components are structural — neither requires informative missingness to operate.

3.7 Permutation calibration and LMM sanity check (EXP-006)

LMM sanity check. Under a single-class homogeneous DGP with no trajectory heterogeneity, R/lme4 returned a Type I error rate of 7.5% (15/200 rejections at $\alpha = 0.05$), within the expected simulation variability around the nominal 5% (95% CI: 4.3–12.2%). Python/statsmodels returned 13.5% (27/200), exceeding the expected range and suggesting anti-conservative behavior in the statsmodels MixedLM implementation — likely attributable to differences in restricted maximum likelihood optimization between implementations. Under the same code with the three-class null DGP, R/lme4 returned 22% Type I error. The 7.5% \rightarrow 22% transition confirms that the LMM false positive inflation is a structural consequence of trajectory heterogeneity, not an implementation artifact.

Permutation calibration. Table 3 presents the results of EXP-006.

Table 3: Permutation calibration results (EXP-006, null scenario, $N = 200/\text{arm}$).

Stress condition	N sims	Parametric Type I	Permutation Type I
Jitter ± 2 months	50	16.0%	4.0%
Rater noise SD = 5	50	4.0%	4.0%

For jitter ± 2 months, the parametric Wald test was anti-conservative (16.0%), but full-pipeline permutation restored nominal Type I error (4.0%). This confirms that the inflated parametric p -value arises from the Wald approximation breaking under heavy visit-time irregularity — the LCMM pipeline itself is unbiased, but the standard error estimates are unreliable when visit times deviate substantially from the nominal schedule. The permutation test, by constructing the null distribution empirically from the observed data structure, is immune to this problem.

For rater noise SD = 5, the parametric Type I error was 4.0%, consistent with the nominal 5% level. The 10% rate observed in EXP-005 ($N = 50$ sims per cell) was attributable to simulation noise rather than systematic inflation, as confirmed by both the larger parametric sample and the permutation results.

These results validate Session 006 Decision 2: full-pipeline permutation inference is mandatory for LCMM-Soft under real-world data conditions. The computational cost ($B \times$ full pipeline evaluations per trial) is substantial but feasible with modern computing resources (approximately 35 minutes for $B = 199$ across 100 trials on 14-core Apple Silicon).

4 Discussion

4.1 Summary of findings

Across six simulation experiments totalling approximately 14,600 simulated trials, complemented by an analytical derivation of collider bias, we found convergent evidence that standard ALS trial endpoints carry a quantifiable statistical cost when the true disease process involves latent trajectory heterogeneity.

The *cost of linearity* — the power penalty from assuming homogeneous linear progression when treatment effects are class-specific — corresponds to an approximately four-fold sample size requirement relative to an oracle class-aware analysis. A practical two-stage LCMM pipeline with pseudo-class inference and Rubin’s rules recovers approximately half this gap, achieving a two-fold sample size penalty relative to the oracle while maintaining Type I error control.

ANCOVA on change from baseline inflates treatment effect estimates approximately ten-fold under the class-specific scenario, even under strictly MAR dropout. This inflation is a structural estimand mismatch — the consequence of targeting the survivor average rather than the population average treatment effect. The analytical derivation confirms that the bias is driven by the covariance between class-specific survival and class-specific treatment effect, a form of collider bias well characterized in the causal inference literature [Hernán, 2010] but not previously quantified for ALS-specific parameters.

Class enumeration via information criteria is complicated by treatment-induced class splitting: the treatment itself creates a new trajectory phenotype (treated-slow \neq untreated-slow) that model selection criteria correctly detect. The solution is to perform class enumeration on pooled data without treatment covariates, then estimate treatment effects within classes in the second stage. ICL provides marginal advantages over BIC at small sample sizes through its entropy penalty.

Critically, these findings are robust to realistic data degradation. The LCMM-Soft pipeline maintained 98–100% power across all eleven stress conditions tested in EXP-005, including a combined severe scenario that simultaneously introduced visit jitter, rater noise, excess dropout, and sporadic missingness. Type I error remained controlled at 2–8% for the majority of conditions. By contrast, the standard LMM exhibited inflated false positive rates (10–36%) across all conditions, including clean data, rendering its power estimates uninterpretable under the class-specific DGP.

4.2 Relationship to existing literature

We emphasize that the individual components of this work are not novel in isolation. Nonlinear ALS progression has been documented [Gordon et al., 2010, van Eijk et al., 2025]. Trajectory heterogeneity has been modelled [Gomeni et al., 2014]. Estimand concerns have been formalized [ICH E9(R1), 2019]. Collider bias from conditioning on post-treatment variables is well understood [Hernán, 2010, Aalen et al., 2015]. Pseudo-class inference has been developed [Vermunt, 2010, Bolck et al., 2004]. LCMMs are mature statistical tools [Proust-Lima et al., 2017].

What has been absent is the *integration*: quantifying the cost of ignoring trajectory heterogeneity in ALS-specific parameters, constructing a complete analysis pipeline with formal Type I error control, and proving — both analytically and via simulation — that ANCOVA bias under heterogeneous progression is structural. Each component builds on established methodology; the contribution is the synthesis and its application to a disease area where the stakes (a $> 97\%$ trial failure rate) are unusually high.

The robustness findings of EXP-005 address a recurring concern in the simulation literature — that novel methods often show advantages only under the idealized conditions of their own data-generating process [Boulesteix et al., 2017]. By demonstrating that the LCMM-Soft pipeline maintains its performance under eleven distinct degradation conditions, including compounding worst-case scenarios, we provide evidence that the method’s advantages are not an artifact of simulation cleanliness. The cross-validation across independent Python and R implementations further strengthens the computational reproducibility of these findings.

We are particularly indebted to the PRECISION-ALS investigators [van Eijk et al., 2025] for providing the most comprehensive evidence to date of trajectory heterogeneity in ALS, and to Proust-Lima et al. [2017] for the LCMM framework that underpins our proposed pipeline.

4.3 Implications for trial design

The implications depend on which treatment effect scenario obtains in reality — which is, of course, unknown at the design stage.

If treatment effects are uniform across trajectory classes, standard LMM analysis is efficient and the LCMM pipeline offers no advantage (indeed, it loses power by splitting the sample). The Holm co-primary strategy mitigates this risk by preserving the overall LMM test.

If treatment effects are class-specific, current ALS trials are operating at approximately one-quarter of their potential efficiency. A trial designed for 80% power under the class-specific scenario with LMM analysis ($N \approx 400$ per arm) would achieve the same power with $N \approx 100$ per arm using an oracle analysis. The practical LCMM pipeline splits the difference: approximately $N \approx 200$ per arm for 80% power.

ANCOVA on change from baseline should not be used as a primary endpoint when survival differs between trajectory classes — which is virtually certain in ALS. At minimum, trialists should explicitly state the estimand targeted by their analysis and acknowledge the survivor average interpretation. The treatment policy estimand, as defined by ICH E9(R1), requires methods that use all available data (such as LMM under MAR, or joint longitudinal-survival models) rather than restricting to survivors.

Pre-specification is essential. The two-stage LCMM pipeline involves numerous analyst degrees of freedom: choice of K_{\max} , information criterion, quality filter thresholds, number of pseudo-class draws, and permutation test parameters. All must be pre-specified to avoid post hoc optimization. Our pre-registration with timestamped commits demonstrates that this is feasible.

4.4 Limitations

Several limitations warrant emphasis.

Simulation only. All results are conditional on the assumed DGP. The three-class structure, trajectory parameters, and survival probabilities are informed by published literature but have not been validated against individual patient-level data. Empirical validation on the PRO-ACT database is planned and will be reported separately. We note, however, that EXP-005 demonstrates the pipeline's robustness extends well beyond the clean DGP specification: performance is maintained under substantial data degradation including visit jitter, rater noise, excess dropout, and missing data. The pipeline's advantages are therefore unlikely to be an artifact of idealized simulation conditions, though empirical validation remains essential.

DGP specificity. The magnitude of the cost of linearity is parameter-dependent. Different class proportions, survival differentials, or treatment effect distributions would yield different power penalties and bias magnitudes. Our results should be interpreted as demonstrating the *existence and mechanism* of these costs under plausible parameters, not as precise predictions for specific future trials.

Simplified LCMM. Our EM-based LCMM implementation uses linear trajectories within classes without link functions or random quadratic terms. Full implementations (e.g., the R package `lcmm`; Proust-Lima et al., 2017) offer more flexible specifications that might improve class recovery and downstream power.

Three-class DGP. The true number of ALS trajectory phenotypes is unknown and may be higher than three, lower than three, or not well represented by discrete classes at all. Continuous heterogeneity (e.g., via random effects with non-normal distributions) is an alternative modelling strategy that our framework does not address.

Permutation test feasibility. Full-pipeline permutation testing with $B = 999$ requires re-

running the complete analysis 999 times per trial. This is computationally intensive and may be impractical for complex LCMM specifications. Parametric bootstrap offers a validated alternative at lower computational cost, though its Type I error properties require separate verification.

Stress-test boundary conditions. Two of the eleven EXP-005 stress conditions initially produced elevated Type I error in the LCMM-Soft pipeline: jitter ± 2 months (16%) and rater noise SD = 5 (10%). EXP-006 subsequently demonstrated that full-pipeline permutation testing resolves both: the jitter outlier was reduced from 16% to 4%, and the rater noise elevation was confirmed as simulation noise. This validates the recommendation that permutation inference should be used for all confirmatory LCMM-Soft analyses, particularly under data conditions that deviate substantially from balanced, regular visit schedules.

AI authorship. This manuscript was produced by an AI research agent. While all analytical decisions were pre-registered and subjected to adversarial deliberation, the agent’s reasoning is ultimately constrained by its training data and architecture. The adversarial deliberation process — in which AI agents with distinct methodological perspectives challenged each other’s reasoning — provides some mitigation but does not substitute for human expert peer review. We welcome scrutiny of both the methods and the process.

4.5 Future work

Empirical validation. Application of the LCMM pipeline to the PRO-ACT database will provide a Trajectory Atlas of ALS progression phenotypes and an empirical test of the class separability assumptions underlying the simulation results. A pre-registered kill switch (median classification entropy < 0.70) will trigger a pivot to a myth-busting paper documenting the limits of latent class approaches.

Contour plots. Power as a function of class separability \times responsive subgroup proportion will characterize the conditions under which the two-stage LCMM pipeline offers meaningful advantages over standard methods, providing a practical tool for trial design.

Adaptive and platform trials. The class enumeration and treatment testing framework is potentially compatible with adaptive designs that update class structure estimates at interim analyses. This extension requires careful consideration of multiplicity and information leakage.

5 Pre-Registered PRO-ACT Analysis Plan

The simulation results presented above motivate but do not validate the proposed pipeline on real patient data. Application to the PRO-ACT database [Atassi et al., 2014] — the largest publicly available repository of ALS clinical trial data — is planned and will be reported separately. The complete analysis protocol was pre-registered via timestamped GitHub commits (commit 75e9221, amended 0b38f6c) prior to any access to patient-level data. PRO-ACT data access was applied for on February 15, 2026 and is pending approval. The locked protocol is summarized here as Section 5 of this preprint, per Session 006 Decision 3.

5.1 Protocol summary

1. **Data harmonization.** Harmonize ALSFRS-R scores and visit dates across trials/eras within PRO-ACT. Define time-zero as randomization/enrollment date. Pre-specified sensitivity analyses: re-align to symptom onset and diagnosis date.

2. **Characterize missingness and death.** Document patterns of missing data and mortality across the dataset. Pre-specify the primary estimand as treatment policy (ICH E9(R1)), with death as an intercurrent event handled via a joint longitudinal-survival model.
3. **Treatment-blind LCMM class enumeration.** Fit LCMMs on pooled data (both treatment arms combined) without treatment as a covariate, for $K = 1$ to $K_{\max} = 5$. Select K by ICL. Minimum class proportion $> 5\%$. Report entropy and posterior classification distributions. Stratify or adjust for trial/era/site to prevent “administrative classes” driven by trial composition rather than biology.
4. **Flexibility on K .** If $K = 2$ rather than 3, proceed. The scientific message is that heterogeneity matters for trial analysis, not that a specific number of classes is correct.
5. **Kill switch.** If median posterior classification entropy < 0.70 , declare that PRO-ACT does not support stable discrete trajectory phenotypes under these endpoints and visit structures. Pivot to pre-specified continuous alternatives: (a) hierarchical random-slope model estimating the treatment effect on the slope distribution using joint-model latent decline parameters; (b) pre-specified quantile treatment effects on individual slopes; (c) permutation-calibrated inference where feasible. Code for this alternative path will be version-controlled and hashed before data access.
6. **Trajectory Atlas.** Construct a standardized figure-and-table set: class trajectories with uncertainty bands, class proportions by trial/era, baseline covariate enrichment (descriptive, not causal), and survival overlays.
7. **Within-class treatment effects.** Estimate treatment effects using LCMM-Soft with pseudo-class draws ($M = 20$) and full-pipeline permutation-calibrated inference ($B = 999$), stratified by trial/study to respect original randomization structures. Clearly separated from the class discovery step.

6 Conclusion

The assumption of linear, homogeneous progression in ALS costs clinical trials statistical power — approximately a four-fold sample size penalty for class-specific treatment effects — and introduces structural bias in common endpoints — approximately ten-fold inflation in ANCOVA estimates from estimand mismatch. These costs are not artifacts of poor data quality or informative missingness; they are mathematical consequences of applying linear models to a nonlinear, heterogeneous disease process and of conditioning on post-treatment colliders.

A pre-specified two-stage LCMM pipeline with ICL-based class enumeration on pooled data, pseudo-class inference via [Vermunt's \(2010\)](#) method, Rubin's variance combination rules, and full-pipeline permutation testing provides a viable alternative with formal Type I error control at a manageable approximately two-fold sample size cost relative to an oracle. Hard class assignment inflates Type I error and should not be used for confirmatory analysis. Stress testing across eleven data degradation conditions — from scheduling jitter and rater noise to compounding worst-case scenarios — confirms that the pipeline's advantages are robust to the messy realities of multi-site clinical trials, with power maintained at 98–100% and Type I error controlled at 2–8% for the majority of conditions.

These findings do not invalidate existing ALS trial results, nor do they claim that methodology alone explains the disease's extraordinary treatment failure rate. They demonstrate that a quantifiable methodological cost exists, that it is largest precisely when the treatment effect is biologically

plausible (concentrated in a responsive subpopulation), and that practical alternatives are available — and robust.

All simulation code, pre-registration records, adversarial deliberation transcripts, and this manuscript are openly available at [GitHub repository URL]. Empirical validation on the PRO-ACT database is forthcoming.

References

- Odd O. Aalen, Richard J. Cook, and Kjetil Røysland. Does Cox analysis of a randomized survival study yield a causal treatment effect? *Lifetime Data Analysis*, 21(4):579–593, 2015.
- Nazem Atassi, James Berry, Amy Shui, Neta Zach, Alexander Sherman, Ervin Sinani, et al. The PRO-ACT database: Design, initial analyses, and predictive features. *Neurology*, 83(19):1719–1725, 2014.
- Christophe Biernacki, Gilles Celeux, and Gérard Govaert. Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(7):719–725, 2000.
- Annabel Bolck, Marcel Croon, and Jacques Hagenaars. Estimating latent structure models with categorical variables: One-step versus three-step estimators. *Political Analysis*, 12(1):3–27, 2004.
- Anne-Laure Boulesteix, Rory Wilson, and Alexander Hapfelmeier. Towards evidence-based computational statistics: Lessons from clinical research on the role and design of real-data benchmark studies. *BMC Medical Research Methodology*, 17(1):138, 2017.
- Jesse M. Cedarbaum, Nancy Stambler, Errol Malta, Cynthia Fuller, Dana Hilt, Bonnie Thurmond, and Atsushi Nakamishi. The ALSFRS-R: A revised ALS functional rating scale that incorporates assessments of respiratory function. *Journal of the Neurological Sciences*, 169(1–2):13–21, 1999.
- Roberto Gomeni, Maurizio Fava, and Zubin Bhagwagar. A latent variable-based approach for determining disease progression in ALS. *Amyotrophic Lateral Sclerosis and Frontotemporal Degeneration*, 15(3–4):276–284, 2014.
- Paul H. Gordon, Bin Cheng, François Salachas, Pierre-François Pradat, Gaëlle Bruneteau, Philippe Corcia, Lucette Lacomblez, and Vincent Meininger. Progression in ALS is not linear but is curvilinear. *Journal of Neurology*, 257(10):1713–1717, 2010.
- Miguel A. Hernán. The hazards of hazard ratios. *Epidemiology*, 21(1):13–15, 2010.
- Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70, 1979.
- ICH E9(R1). Addendum on estimands and sensitivity analysis in clinical trials to the guideline on statistical principles for clinical trials. Technical report, International Council for Harmonisation, 2019.
- Dmitry Petrov, Carol Mansfield, Alain Moussy, and Olivier Hermine. ALS clinical trials review: 20 years of failure. Are we any closer to registering a new treatment? *Frontiers in Aging Neuroscience*, 9:68, 2017.

Cécile Proust-Lima, Viviane Philipps, and Benoit Liquet. Estimation of extended mixed models using latent classes and latent processes: The R package `lcmm`. *Journal of Statistical Software*, 78(2):1–56, 2017.

Donald B. Rubin. *Multiple Imputation for Nonresponse in Surveys*. Wiley, 1987.

Ruben P. A. van Eijk et al. PRECISION-ALS: Clustering of ALSFRS-R trajectories reveals heterogeneous disease progression. *Neurology*, 2025. Forthcoming.

Jeroen K. Vermunt. Latent class modeling with covariates: Two improved three-step approaches. *Political Analysis*, 18(4):450–469, 2010.

A Complete Simulation Code

Complete simulation code (Python and R). Available at [GitHub repository URL].

B Analytical Derivation of K-Class Collider Bias

Full analytical derivation of K -class collider bias, including general conditions for unbiasedness and extensions to joint longitudinal-survival models.

C Pre-Registration Records

Pre-registration records with timestamped GitHub commits, including amendments with rationale.

D Adversarial Deliberation Transcripts

Adversarial deliberation transcripts (Board Room Sessions 001–006). These document the structured process by which analytical decisions were debated, challenged, and refined by multiple AI agents with distinct methodological perspectives.

E Extended Simulation Results

All scenarios, sample sizes, and methods; K -selection distributions; MNAR gradient details.

F EXP-005 Stress-Test Detailed Results

Per-condition power, Type I error, point estimates, and standard errors for both LCMM-Soft and LMM across all eleven degradation conditions, with Python–R cross-validation comparison.

G EXP-006 Permutation Calibration Results

EXP-006 permutation calibration results and LMM sanity check code (R/lme4).

Manuscript prepared 2026-02-17. Revised 2026-02-17 (v3: added EXP-006 permutation calibration, LMM sanity check, pre-registered PRO-ACT protocol). Pre-registration: GitHub commit 75e9221. All code open source.