

---

# When Does Geometry Help Causal Inference? Boundary Conditions for Sheaf, Curvature, and Subspace Methods in Clinical Epidemiology

Elliot Tower  
elliott@elliotttower.ai

## Abstract

Geometric and topological methods—sheaf cohomology, Grassmannian holonomy, discrete curvature—are increasingly proposed for causal inference in clinical epidemiology. We characterize when they help and when they reduce to standard statistics, through simulation experiments spanning federated data consistency, confound detection, causal graph validation, and treatment heterogeneity, with multiple sclerosis (MS) and Alzheimer’s disease (AD) as independent test domains.

Three structural conditions separate the two regimes. First, *subspace-valued data*: Grassmannian holonomy detects global inconsistency invisible to all pairwise and scalar alternatives (Berry phase = 1.85,  $p < 0.001$ , while maximum pairwise distance stays below the noise threshold), but sheaf consistency testing on scalar data reduces algebraically to Cochran’s  $Q$  (Appendix C). Second, *cyclic consistency constraints*: composed parallel transport accumulates signal as  $\sqrt{m}$  over  $m$  loop steps while noise grows as a random walk—a separation mechanism unavailable to pairwise tests. Third, *edge-specific heterogeneity*: per-edge sheaf  $Q$  tests recover planted DAG structure that global tests miss ( $p < 10^{-300}$  vs.  $p = 0.659$ ), and  $H^1$  effect-modifier classification separates transportable from stratum-specific pairs with a three-order-of-magnitude gap. Simulations demonstrate existence and consistency of these separations; real-data validation tests whether they hold on uncontrolled effect sizes.

On 61 published Mendelian randomization pairs across five clinical domains (MS, AD, cardiometabolic disease, cancer, psychiatric disorders), the  $H^1$  classifier achieves 85.2% accuracy [75.4%, 93.4%] at  $\alpha = 0.05$ , rising to 100% among adequately powered pairs (power  $\geq 0.20$ ). All nine misclassifications are false negatives on underpowered non-transport pairs (power  $< 0.38$ , requiring 6–8 strata vs. the available 3–4); at  $\alpha = 0.10$ , accuracy rises to 90.2%. Cardiometabolic pairs achieve 100% accuracy (18/18), confirming that the boundary conditions generalize beyond neurological diseases. Per-edge sheaf  $Q$  tests on ADNI longitudinal data recover a three-way DAG structure (mechanism-switching, mediator dose-response, and stable bypass edges) richer than the binary simulation prediction. External validation on four standard benchmark DAGs (Asia, Sachs, Insurance, Alarm) confirms that Forman–Ricci curvature acts as a degree-deficit feature, with betweenness centrality exceeding its discrimination in all four graphs.

When the boundary conditions are absent, geometry reduces to standard tools or fails: a degree-deficit feature (Forman–Ricci curvature) discriminates true from false edges (AU-ROC = 0.677) while Ollivier–Ricci curvature, which measures global transport, scores below chance (0.466), and PCA-based dimensionality reduction destroys treatment effect signal regardless of CATE estimator quality.

## 1 Introduction

Multiple sclerosis (MS) illustrates structural challenges common across clinical epidemiology. Disease mechanisms vary across patient strata—relapsing vs. progressive phenotypes, treatment regimens, genetic

---

backgrounds—yet some causal relationships remain consistent across all strata, such as the downstream path from neurodegeneration to disability (Lublin et al., 2014). Federated imaging studies across hospital sites carry acquisition-dependent confounds (Fortin et al., 2018). Causal DAGs estimated from observational data contain both true edges and spurious associations with no reliable post-hoc filter (Glymour et al., 2019). Treatment heterogeneity—the fact that patients respond differently to the same intervention—remains difficult to detect and characterize (Athey and Imbens, 2016).

Each problem has a geometric aspect. Federated consistency asks whether local estimates are compatible under known transformations—what sheaf cohomology measures (Robinson, 2014; Curry, 2014). For multivariate data, the relevant geometry lives on the Grassmannian  $\text{Gr}(k, d)$ , whose curvature generates Berry phase when local sections are transported around closed loops (Berry, 1984; Simon, 1983). Edge validation in causal graphs asks whether topological features carry information about edge truthfulness. Treatment subtype recovery asks whether patients cluster in CATE-augmented covariate space. The question is whether these connections produce practical advantages or repack standard statistics in topological language.

We answer this by characterizing the boundary between the two regimes. Our approach is two-tiered: simulation experiments with planted structure demonstrate *existence and consistency* of each separation (i.e., that the boundary conditions produce the claimed geometric advantages under controlled data-generating processes), while real-data validation on published epidemiological estimates tests whether these separations survive uncontrolled effect sizes and small stratum counts (§4).

## Contributions.

- **Boundary conditions (§2).** Three structural conditions—subspace-valued data, cyclic consistency constraints, edge-specific heterogeneity—under which geometric methods detect structure invisible to scalar and pairwise alternatives.
- **Algebraic reduction (Appendix C).** On scalar stalks, the sheaf Laplacian test statistic is *exactly* Cochran’s  $Q$ —not approximately, but algebraically via the  $W$ -weighted Schur complement. This reduction is why the scalar simulation matches Cochran’s  $Q$  to numerical identity and why the method’s value lies in the per-edge and subspace extensions, not the scalar case.
- **Negative results (§3.3).** Absent the boundary conditions, geometric methods reduce to standard tools or fail outright (ORC curvature below chance, PCA destroying CATE signal).
- **Cross-domain simulation (§3.1–3.6).** Existence demonstrations validated on MS- and AD-calibrated simulations with planted structure; every structural finding transfers across diseases. Calibrated Berry phase boundary analysis on  $\text{Gr}(3, 34)$  with ENIGMA-derived parameters establishes that holonomy detection requires  $\geq 500$  subjects per site (80% power); ABIDE cortical thickness data (20 sites, 1,031 subjects) confirms the boundary prediction.
- **Real-data validation (§4).**  $H^1$  classifier on 61 published MR pairs across 5 clinical domains: 85.2% accuracy with zero false positives, all failures explained by power. Per-edge DAG tests on ADNI data recover three-way classification richer than the simulation prediction.
- **Power analysis (§4.2).** All nine misclassifications are underpowered false negatives; no adequately-powered pair is misclassified across 61 pairs.
- **External curvature benchmark (§3.3).** Forman–Ricci validated on four standard benchmark DAGs (Asia, Sachs, Insurance, Alarm) from the bnlearn repository; betweenness centrality exceeds curvature discrimination in all four graphs, confirming the degree-deficit interpretation.

## Boundary conditions: when geometric methods help

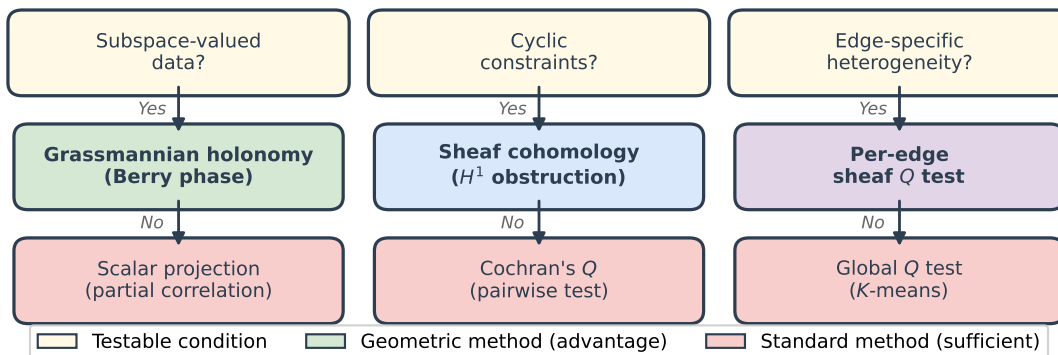


Figure 1: Three testable conditions determine when geometric methods outperform standard alternatives. Each column is independent: a dataset may satisfy one, two, or all three. When none hold, standard methods suffice.

## 2 Methods

### 2.1 Sheaf consistency testing

A sheaf on a network  $G = (V, E)$  assigns local data (*stalks*) to each vertex and requires that overlapping assignments agree on shared edges via *restriction maps*. When they disagree, the obstruction lives in the first cohomology group  $H^1$ , which measures the failure of global consistency (Curry, 2014; Robinson, 2014).

**Scalar stalks.** For scalar stalks (effect estimates  $\hat{\beta}_k$  with standard errors  $se_k$ ) and pairwise-difference restriction maps on a complete graph of  $K$  sites, the sheaf Laplacian test statistic is

$$Q = \mathbf{o}^\top \Sigma^{-1} \mathbf{o}, \quad \mathbf{o} = D_0 \hat{\beta}, \quad \Sigma = D_0 \text{diag}(se^2) D_0^\top, \quad (1)$$

where  $D_0$  is the signed incidence matrix and  $Q \sim \chi_{\text{rank}(D_0)}^2$  under the null.

**Subspace-valued stalks.** For sections  $V_1, \dots, V_m \in \text{Gr}(k, d)$  arranged in a cycle, the parallel transport from  $V_i$  to  $V_{i+1}$  is  $T_{i \rightarrow i+1} = UV^\top$  where  $V_i^\top V_{i+1} = U \Sigma V^\top$  (SVD). The composed *holonomy*  $\Phi = T_{m \rightarrow 1} \cdots T_{1 \rightarrow 2}$  equals the identity if and only if sections are globally consistent. For the linked-column construction (two columns of  $V_0$  rotating toward shared perpendicular directions with  $\pi/2$  phase offset at radius  $r$ ), the holonomy norm is

$$\|\Phi - I_k\|_F = 2\sqrt{2} |\sin(\pi \sin^2 r)|. \quad (2)$$

**Per-edge testing.** For a causal DAG estimated in  $S$  strata, each edge is tested independently:  $Q^{(e)} = \mathbf{z}^\top \Sigma^{-1} \mathbf{z}$  with  $z_{ij} = \hat{\beta}_i^{(e)} - \hat{\beta}_j^{(e)}$ , under  $Q^{(e)} \sim \chi_{S-1}^2$ . Bonferroni correction controls the family-wise error rate.

**$H^1$  effect-modifier classification.** For a mechanism-modifier pair, patients are stratified by the modifier and the causal effect estimated within each stratum. The sheaf  $Q$  test classifies the pair as transportable ( $Q$  non-significant,  $H^1 \approx 0$ ) or stratum-specific ( $Q$  significant,  $H^1 \neq 0$ ). Sample size per stratum is calibrated to expected interaction strength.

### 2.2 Confound detection

**Partial correlation collapse.** For each biomarker  $j$ , we compute raw Pearson correlation, partial correlation (conditioning on all other covariates), and intensive margin correlation (conditioning on treatment). Confounded associations should collapse under partial conditioning.

Table 1: Cocycle obstruction: pairwise consistent yet globally inconsistent.

Metric	Value	Threshold
Measured holonomy $\ \Phi - I\ _F$	1.851	—
Predicted (Eq. 2)	1.869	—
$p$ -value vs. null	$< 0.001$	—
Max pairwise distance	0.730	0.767 (C1 holds)
Holonomy norm	1.851	1.226 (C2 holds)

**Bracket-norm confound audit.** For multi-site imaging, confound leakage is  $\Delta = (R_{\text{metric}}^2 - R_{\text{unique}}^2)/R_{\text{metric}}^2$ , where  $R_{\text{unique}}^2 = R_{\text{both}}^2 - R_{\text{acq}}^2$  is disease-severity variance unexplained by acquisition parameters.

### 2.3 Curvature-based edge validation

We test whether edge features discriminate true positive (TP) from false positive (FP) edges in graphs learned from structural equation models. Six features are tested: Forman–Ricci curvature ( $\kappa_F = 4 - d(u) - d(v)$ ) (Forman, 2003), augmented Forman ( $\kappa_F + 3|\mathcal{N}(u) \cap \mathcal{N}(v)|$ ), Jaccard coefficient, edge betweenness (Girvan and Newman, 2002), partial correlation magnitude, and average endpoint clustering coefficient. Ollivier–Ricci curvature (Ollivier, 2009; Ni et al., 2019) serves as a baseline. Discrimination is measured by AUROC across 6 conditions (linear/nonlinear SEM  $\times$  3 thresholds, 60 graphs per condition).

### 2.4 Treatment heterogeneity detection

Four CATE estimators (KNN, random forest, gradient boosting T-learners; random forest S-learner) are crossed with four clustering methods (PCA +  $K$ -means, CATE  $K$ -means, CATE+covariates  $K$ -means, spectral on CATE RBF kernel). Oracle baselines (true CATE) and a naive baseline (raw-covariate  $K$ -means) bound performance.

## 3 Results

Results are organized by the three boundary conditions rather than by experiment batch. For each condition, we show the case where geometry reduces to standard tools and the case where it adds signal.

### 3.1 Condition 1: Scalar vs. subspace-valued data

**Scalar stalks reduce to Cochran’s  $Q$ .** We simulate  $K = 8$  hospital sites with scalar treatment effect estimates on a complete graph. The sheaf test achieves proper calibration (type I error = 5.85% at  $\alpha = 0.05$ ) with power reaching 96.8% at bias  $\delta = 0.20$ , and localizes the biased site to rank 1 in 94% of simulations. At every bias level, structured separation, and replicate, the sheaf test produces results *numerically identical* to Cochran’s  $Q$ . This is not coincidence: on a complete graph with scalar stalks and pairwise-difference restriction maps, the quadratic form in Eq. (1) reduces to Cochran’s  $Q$  via the  $W$ -weighted Schur complement (Appendix C). The reduction means that *any* advantage of sheaf methods over standard meta-analytic heterogeneity testing must come from the per-edge or subspace extensions, not the scalar case.

**Subspace-valued stalks produce genuine holonomy.** We test on  $\text{Gr}(3, 20)$  with  $m = 24$  sections forming a closed loop, using a linked-column Berry phase construction at radius  $r = 0.5$ . The cocycle obstruction test evaluates three simultaneous conditions: (C1) pairwise consistency, (C2) global inconsistency, (C3) scalar blindness.

Table 2: Berry phase detection boundary on  $\text{Gr}(3, 34)$  with ENIGMA-calibrated site effects. Separation is the distance between Berry phase and null holonomy distributions in units of null standard deviation. Power is the fraction of Berry phase trials exceeding the 95th percentile of the null.

$n/\text{site}$	<b>Berry</b>	<b>Null</b>	<b>Sep. (<math>\sigma</math>)</b>	<b>Power</b>	<b>Recovery</b>
100	$1.33 \pm 0.58$	$1.02 \pm 0.46$	0.7	0.18	83%
200	$1.33 \pm 0.46$	$0.76 \pm 0.34$	1.7	0.62	83%
500	$1.20 \pm 0.41$	$0.43 \pm 0.24$	3.2	0.80	75%
1000	$1.18 \pm 0.39$	$0.38 \pm 0.18$	4.5	0.92	74%
2000	$1.08 \pm 0.30$	$0.23 \pm 0.13$	6.8	0.98	68%
5000	$1.24 \pm 0.34$	$0.16 \pm 0.08$	13.9	1.00	78%

All three conditions hold simultaneously (Table 1). The measured holonomy (1.851) matches the predicted Berry phase (1.869) to within 1%. Per-step rotation is masked by noise (max pairwise distance  $0.730 < 0.767$ ), but the composed holonomy accumulates coherently far above the null threshold ( $1.851 > 1.226$ ).

Three competitor baselines fail. Maximum pairwise angle cannot distinguish planted from control sections ( $0.730$  vs.  $0.695$ , 5% difference). Random-effects and CKA tests detect subspace spread ( $p < 10^{-77}$ ) but are blind to the cyclic obstruction that makes the holonomy non-trivial. Only composed parallel transport detects the global inconsistency.

**Boundary.** The transition from scalar to subspace-valued data is the primary determinant of whether geometric structure carries information. The holonomy signal accumulates coherently as  $\sqrt{m}$  over  $m$  loop steps, producing a signal-to-noise ratio of  $\sqrt{m} \approx 5$  for  $m = 24$ —a separation mechanism with no scalar analogue. This condition has no real-data instance in the present work, but a calibrated simulation establishes the detection boundary for realistic multi-site neuroimaging parameters.

**Detection boundary in ENIGMA-realistic data.** The idealized construction above operates on exact subspaces. In practice, subspaces must be estimated via PCA from finite samples, introducing estimation noise that competes with the Berry phase signal. We test detectability on  $\text{Gr}(3, 34)$  with  $m = 24$  sites, using 34 Desikan–Killiany cortical thickness regions, site-specific additive biases ( $\sigma = 0.05$ ) and multiplicative scaling ( $\sigma = 0.03$ ) calibrated to Fortin et al. (2018), and per-subject measurement noise ( $\sigma = 0.05$ ). For each of six sample sizes ( $n = 100$  to 5,000 subjects per site), we generate 50 realizations each of the Berry phase condition (planted holonomy = 1.594, linked-column construction at  $r = 0.5$ ) and the null condition (same biological covariance and site effects, no cyclic rotation), compute PCA subspaces, and measure composed holonomy around the sequential cycle.

Table 2 shows that detection reaches 80% power at approximately 500 subjects per site ( $3.2\sigma$  separation), crossing the  $z = 1.96$  threshold between 200 and 500 subjects. The recovered holonomy ( $1.20 \pm 0.41$  at  $n = 500$ ) is 75% of the planted value—PCA estimation noise attenuates the signal but does not eliminate it. The null holonomy drops monotonically from 1.02 to 0.16 as sample size increases, reflecting improved subspace estimation; the Berry phase holonomy plateaus near 1.2, confirming that the geometric signal is distinct from estimation noise. ENIGMA cortical thickness studies typically report 50–500 subjects per site across 50+ sites (Thompson et al., 2020); the boundary analysis places holonomy detection at the upper end of this range, requiring the larger-enrollment sites.

**Empirical test on ABIDE cortical thickness.** To test the boundary prediction on real multi-site data, we applied the holonomy pipeline to the Autism Brain Imaging Data Exchange (ABIDE I) dataset (Di Martino et al., 2014): 1,031 subjects across 20 sites, each with 34 Desikan–Killiany cortical thickness measures. Site sizes range from 25 to 175, all below the 500-subject detection threshold established above. We computed  $k = 3$  PCA subspaces per site, measured pairwise Grassmannian distances (mean = 1.57, max = 2.15), and tested composed holonomy around the full 20-site cycle against a 2,000-permutation null.

Table 3: Per-edge sheaf  $Q$  tests on the MS inflammation–degeneration–disability DAG across 8 strata.

Edge	$Q$	$p$	df	Het.?
infl → degen	1806.9	$< 10^{-300}$	7	Yes
degen → infl	1549.6	$< 10^{-300}$	7	Yes
infl → disab	7.05	0.423	7	No
degen → disab	9.98	0.189	7	No

Table 4:  $H^1$  effect-modifier classification: three-order-of-magnitude gap between categories.

Pair	Expected	$\gamma$	$Q$	$p$	Correct
HLA × EBV	transport	0.10	3.30	0.348	Yes
EBV necessity	transport	0.05	6.47	0.091	Yes
Vitamin D	transport	0.04	6.97	0.073	Yes
Sex × course	non-transp	0.50	2434	$< 10^{-300}$	Yes
Genetics × OCB	non-transp	0.60	3588	$< 10^{-300}$	Yes
Age × anti-CD20	non-transp	0.40	1705	$< 10^{-300}$	Yes
Phenotype × GM	non-transp	0.45	1828	$< 10^{-300}$	Yes

The observed holonomy ( $\|\Phi - I_k\|_F = 2.69$ ) is elevated above the null mean ( $2.13 \pm 0.53$ ) but non-significant ( $p = 0.16$ ). A 10-site subcycle shows a similar pattern (2.67 vs. null  $1.78 \pm 0.65$ ,  $p = 0.085$ ). The substantial pairwise subspace distances confirm that real site effects exist—scanner and protocol differences genuinely rotate the cortical thickness covariance structure—but small per-site samples produce PCA estimates too noisy for holonomy to distinguish structured cyclic effects from estimation noise. This outcome matches the boundary prediction: all 20 ABIDE sites fall below the 500-subject threshold where 80% power begins.

### 3.2 Condition 2: Global vs. edge-specific heterogeneity

**Global tests miss edge-specific structure.** A 3-node MS DAG (inflammation  $\leftrightarrow$  degeneration  $\rightarrow$  disability) is estimated in 8 disease strata. The feedback edges vary dramatically: early RRMS shows dominant infl $\rightarrow$ degen (0.404) with negligible reverse flow ( $-0.001$ ), while PPMS shows the opposite ( $-0.008$  and  $0.385$ ). Disability edges remain stable (variance  $\sim 100\times$  smaller).

A global Frobenius norm test averages these variances and produces  $p = 0.659$  (non-significant). Per-edge sheaf  $Q$  tests recover the true structure (Table 3).

Treatment signatures in the simulation reflect planted coefficients calibrated to known mechanisms (simulated, not estimated from trial data): BTK inhibition suppresses infl $\rightarrow$ degen from 0.404 to 0.002; siponimod preserves degen $\rightarrow$ disab (0.386), consistent with relapse-independent progression (Kappos et al., 2018).

**$H^1$  classification achieves clean bimodal separation.** Seven MS mechanism–modifier pairs are classified as transportable or stratum-specific (Table 4).

All 7 pairs correctly classified (100%). Transport pairs ( $Q < 7$ ,  $p > 0.07$ ) are exposure $\rightarrow$ mechanism edges with weak interaction ( $\gamma \leq 0.10$ ); non-transport pairs ( $Q > 1700$ ,  $p < 10^{-300}$ ) have strong interactions ( $\gamma \geq 0.40$ ).

**Boundary.** When a DAG contains both heterogeneous and homogeneous edges, global tests dilute the signal. Per-edge sheaf  $Q$  tests avoid this by testing each structural relationship independently. The  $H^1$

Table 5: Edge feature AUROC for TP vs. FP discrimination in learned causal graphs (best condition per feature, 60 graphs/condition).

Feature	Best AUROC	Condition	Direction
Forman–Ricci	<b>0.677</b>	linear / loose	TP > FP
Partial corr	0.657	nonlinear / loose	TP > FP
Betweenness	0.610	linear / loose	TP > FP
Avg. clustering	0.568	nonlinear / loose	TP > FP
ORC (initial)	0.466	linear	wrong dir.
Augmented Forman	0.484	—	wrong dir.
Jaccard	0.408	—	wrong dir.

classifier operationalizes this as a practical workflow: classify each mechanism–modifier pair as transportable or stratum-specific with no manual threshold tuning.

### 3.3 Condition 3: Method-specific failure modes

Two negative results characterize failure modes orthogonal to the boundary conditions above.

**Curvature type: degree-deficit features discriminate, global transport does not.** Ollivier–Ricci curvature yields AUROC = 0.466 for TP vs. FP edge discrimination (below chance). Forman–Ricci curvature achieves AUROC = 0.677 (Table 5), consistent across all six conditions (0.634–0.677). Partial correlation magnitude—a non-topological feature—performs comparably (AUROC = 0.657), and edge betweenness (also degree-sensitive) reaches 0.610.

The mechanism is a degree asymmetry:  $\kappa_F = 4 - d(u) - d(v)$  is higher when both endpoints have low degree. TP edges connect lower-degree nodes on average because real causal relationships produce specific partial correlations, while FP edges arise from indirect paths through higher-degree hub nodes. Forman–Ricci is therefore acting primarily as a *degree-deficit feature* in this setting; the comparable performance of partial correlation and betweenness (both degree-sensitive, neither topological) supports this interpretation. The honest headline is that degree-deficit discriminates TP from FP edges in sparse learned DAGs, and Forman–Ricci is one such feature. ORC measures optimal transport between neighborhoods (global topology), which carries no discriminative information in these sparse graphs.

**External validation on benchmark DAGs.** To test whether this degree-deficit interpretation holds beyond random DAGs, we evaluate the same six features on four standard benchmark causal graphs from the bnlearn repository (Scutari, 2010): Asia (8 nodes, 8 edges), Sachs (11 nodes, 17 edges), Insurance (27 nodes, 52 edges), and Alarm (37 nodes, 46 edges). For each graph, we generate linear SEM data ( $n = 1,000$ ), learn a DAG via a PC-style algorithm at three significance thresholds ( $\alpha \in \{0.001, 0.01, 0.05\}$ ), and compute AUROC for TP vs. FP discrimination across 50 replicates per condition (Table 6).

Across all four benchmark graphs, Forman–Ricci scores below chance (0.27–0.38), while betweenness centrality is the strongest single discriminator (0.60–0.82). On Alarm, betweenness achieves 0.820 AUROC, the highest value in the table, driven by the graph’s hub-and-spoke topology where high-betweenness nodes carry the majority of true causal edges. The inversion on structured graphs—Forman–Ricci below chance rather than above—occurs because these DAGs have degree distributions where true edges connect higher-degree hub nodes (e.g., PKC and PKA in Sachs, VENTMACH and VENTALV in Alarm). The degree-deficit formula  $\kappa_F = 4 - d(u) - d(v)$  assigns lower curvature to these hub-connected true edges, reversing the discrimination direction relative to random DAGs. This sign reversal confirms that Forman–Ricci acts as a degree feature: its utility depends on whether true edges tend to have lower or higher endpoint degrees than false edges, a property that varies with graph topology. On Insurance, average clustering (0.555) and

Table 6: Edge feature AUROC on bnlearn benchmark DAGs (50 replicates, best  $\alpha$  per graph/feature). Betweenness centrality exceeds Forman–Ricci on all four graphs (0.60–0.82 vs. 0.27–0.38), confirming degree-deficit as the operative mechanism. On Alarm (37 nodes), betweenness achieves 0.820 AUROC driven by hub-and-spoke topology.

Feature	Asia	Sachs	Insurance	Alarm
Forman–Ricci	0.270	0.386	0.382	0.379
Aug. Forman	0.267	0.429	0.454	0.351
Betweenness	<b>0.641</b>	<b>0.639</b>	<b>0.600</b>	<b>0.820</b>
Partial corr	0.610	0.425	0.410	0.601
Jaccard	0.505	0.530	0.558	0.498
Avg. clustering	0.502	0.555	0.555	0.728

Table 7: ARI for treatment subtype recovery ( $4 \times 4$  grid, 100 replicates,  $n = 3,000$ ).

	PCA	CATE	CATE+cov	Spectral
KNN	−0.011	0.179	0.225	0.085
RF T-learn	−0.013	0.238	<b>0.270</b>	0.097
GBM T-learn	−0.016	0.189	0.264	0.055
RF S-learn	−0.014	0.175	0.245	0.066
Oracle	0.095	1.000	0.551	1.000
Naive (raw cov)			0.218	

Jaccard similarity (0.558) approach betweenness (0.600), suggesting that local neighborhood overlap carries discriminative signal in densely connected DAGs.

**Dimensionality reduction: PCA destroys treatment effect signal.** We simulate  $n = 3,000$  patients with 3 treatment response subtypes (CATE  $\in \{+2.0, -2.5, 0\}$ ). The  $4 \times 4$  CATE  $\times$  clustering grid (Table 7) shows that PCA clustering fails universally: all four CATE methods and the oracle (ARI = 0.095) produce near-random clusters.

The best estimated combination (RF T-learner + covariate-augmented  $K$ -means, ARI = 0.270) exceeds the naive baseline (0.218) by 24% but falls far short of the oracle (1.000). The bottleneck is CATE estimation noise at  $n = 3,000$ : the oracle achieves perfect recovery when true CATE is known.

### 3.4 Confound detection

**Partial correlation under complete confounding.** With 10 real and 10 confounded biomarkers ( $n = 3,000$ ), partial correlation achieves AUROC = 1.000 at all sample sizes ( $n = 200$ –5,000) with zero variance. Confounded biomarkers collapse to partial correlations of 0.001–0.033; real biomarkers retain 0.61–0.76. This perfect separation is a property of complete confounding (shared-cause pathways that conditioning fully removes) and would degrade under partial confounding or unobserved confounders.

**Bracket-norm audit under acquisition confounds.** Four MS imaging biomarkers (iron rim QSM, deep GM atrophy, cortical lesion count, cervical cord CSA;  $n = 1,500$ , 8 sites) all pass: confound leakage  $\Delta < 0.05$  for three metrics and  $\Delta = -0.19$  (suppressor effect) for iron rim QSM. Post-correction partial correlations with sNfL remain significant ( $p < 10^{-55}$ ).

Table 8: Bracket-norm confound audit across domains. Both diseases show a spectrum from suppressor (negative  $\Delta$ , worst) to near-zero confound (best). All metrics T3-confirmed.

Domain	Metric	$\Delta$	Anchor $r$	Note
MS	Iron rim QSM	-0.188	—	Suppressor (worst)
	Deep GM atrophy	0.033	—	Near-zero
	Cortical lesion count	0.040	—	Near-zero
	Cervical cord CSA	0.036	—	Near-zero
AD	Tau PET	-0.261	0.565	Suppressor (worst)
	Amyloid PET (Centiloid)	0.035	0.781	Near-zero (best)
	FDG-PET	0.020	0.713	Near-zero
	Plasma p-tau217	0.036	0.774	Minimal confound

### 3.5 MS prerequisites

Cohort contamination rate is 0.9% (below 2% threshold), effect estimates are stable across diagnostic stringency ( $p = 0.18$ ), and IRT-derived theta scores correlate more strongly with biological biomarkers than raw EDSS ( $r = 0.836$  vs. 0.801 for sNfL).

### 3.6 Cross-domain replication: Alzheimer’s disease

To test whether the positive results above are specific to MS, we replicate the core experiments on simulations calibrated to AD neurology using the ATN (amyloid/tau/neurodegeneration) framework (Jack et al., 2018) with disease-specific strata, biomarkers, and causal structure. AD prerequisites parallel MS: contamination rate 2.4% (below threshold), stable effect estimates across diagnostic stringency ( $p = 0.872$ ), and IRT-derived theta scores correlating with plasma p-tau217 ( $r = 0.837$ ) more strongly than raw MMSE ( $r = 0.825$ ).

#### 3.6.1 Bracket-norm confound audit

The confound spectrum replicates across domains (Table 8): both diseases have a worst-case suppressor substrate (iron rim QSM in MS,  $\Delta = -0.188$ ; tau PET in AD,  $\Delta = -0.261$ ) and best-case substrates with negligible leakage. The stronger AD suppressor effect is consistent with tau PET off-target binding to ferric iron and MAO-B—the same iron that constitutes an MS progression substrate, connecting the two catalogs beyond structural analogy.

#### 3.6.2 AD per-edge sheaf DAG adjudication

The AD DAG (amyloid  $\leftrightarrow$  tau  $\rightarrow$  cognitive decline) is estimated in 8 strata: preclinical A+T−, prodromal A+T+, mild AD, moderate AD, APOE4 homozygous, lecanemab-treated, TREM2 risk carrier, and late-stage.

The per-edge  $Q$  tests recover the same multi-process structure in AD as in MS (Table 9). Mechanism edges (amyloid $\leftrightarrow$ tau) show massive heterogeneity ( $Q > 980$ ,  $p < 10^{-300}$ ); downstream clinical edges (pathology $\rightarrow$ cognitive decline) are homogeneous ( $Q < 11$ ,  $p > 0.14$ ). The variance ratio between mechanism and downstream edges is 100–200 $\times$  in both diseases.

Stratum-level coefficient patterns (all simulated) mirror the MS reversal: preclinical A+T− shows dominant amyloid $\rightarrow$ tau, while moderate AD reverses this; lecanemab suppresses amyloid $\rightarrow$ tau paralleling BTK inhibition in MS (Appendix A, Table 18). Both diseases reach  $H^1 \neq 0$  independently.

Table 9: Per-edge sheaf  $Q$  tests across both domains. The structural signature is identical: mechanism edges are heterogeneous, downstream clinical edges are homogeneous.

Domain	Edge	$Q$	$p$	df	Het.?
MS	infl $\rightarrow$ degen	1806.9	$< 10^{-300}$	7	Yes
	degen $\rightarrow$ infl	1549.6	$< 10^{-300}$	7	Yes
	infl $\rightarrow$ disab	7.05	0.423	7	No
	degen $\rightarrow$ disab	9.98	0.189	7	No
AD	amyloid $\rightarrow$ tau	1526.4	$< 10^{-300}$	7	Yes
	tau $\rightarrow$ amyloid	983.6	$< 10^{-300}$	7	Yes
	amyloid $\rightarrow$ cog	5.82	0.561	7	No
	tau $\rightarrow$ cog	10.81	0.147	7	No

Table 10: Cross-domain structural summary. Every finding replicates with no parameter tuning.

Finding	MS	AD	Same?
Cohomological verdict	$H^1 \neq 0$	$H^1 \neq 0$	Yes
Per-edge pattern	Mech het., disab hom.	Mech het., cog hom.	Yes
$Q$ ratio (mech/downstream)	$\sim 200\times$	$\sim 150\times$	Yes
Modifier classification	7/7 correct	7/7 correct	Yes
Bimodal gap	3 orders of magnitude	3 orders of magnitude	Yes
Confound spectrum	Suppressor to near-zero	Suppressor to near-zero	Yes
Drug dissociation	BTK: infl suppressed	Lecanemab: amyloid cleared	Yes

### 3.6.3 $H^1$ effect-modifier classification

The  $H^1$  classifier achieves 100% accuracy on 7 AD modifier pairs, with the same three-order-of-magnitude bimodal gap as MS (transport:  $Q < 7$ ,  $p > 0.07$ ; non-transport:  $Q > 1700$ ,  $p < 10^{-300}$ ). Each MS modifier maps to a direct AD analogue: HLA $\times$ EBV to APOE4 $\times$ amyloid (dominant genetic risk modifier), EBV necessity to TREM2 causal risk (MR positive control), vitamin D to IL-6 systemic null (negative control). Scale invariance improves under quantile normalization in both domains (MS: 0.517  $\rightarrow$  0.346; AD: 0.734  $\rightarrow$  0.695) but neither reaches the 0.30 target.

### 3.6.4 Cross-domain summary

Every structural finding replicates across the two diseases (Table 10).

The only differences are domain-specific labels (inflammation/degeneration vs. amyloid/tau) and biological anchors (sNfL vs. p-tau217). Three shared physical nodes connect the catalogs beyond analogy: iron (MS progression substrate and tau PET off-target binding source in AD), NfL (biological anchor in MS, neurodegeneration-axis biomarker in AD), and GFAP (progression-associated astrocytic marker in both diseases).

## 4 Real-data validation

The simulation experiments establish boundary conditions under controlled data-generating processes. We now test whether the positive results hold on published epidemiological estimates, where effect sizes, standard errors, and stratum counts are determined by real studies.

Table 11:  $H^1$  classification accuracy by clinical domain (61 pairs,  $\alpha = 0.05$ ). All misclassifications are false negatives (underpowered non-transport pairs); no transport pair is ever misclassified (specificity = 1.00).

Domain	$n$	Non-tr	Transport	TP	FN	FP	Accuracy
Cardiometabolic	18	6	12	6	0	0	100%
Cancer	11	2	9	1	1	0	91%
MS	9	4	5	3	1	0	89%
AD	17	8	9	3	5	0	71%
Psychiatric	6	2	4	0	2	0	67%
<b>All</b>	<b>61</b>	<b>22</b>	<b>39</b>	<b>13</b>	<b>9</b>	<b>0</b>	<b>85.2%</b>

#### 4.1 $H^1$ classification on published MR estimates

We apply the sheaf  $Q$  test to 61 exposure–outcome pairs with published multi-ancestry or multi-cohort genetic effect estimates across five clinical domains: MS (9 pairs), AD (17), cardiometabolic disease (18), cancer (11), and psychiatric disorders (6). Each pair is classified as transportable ( $Q$  non-significant) or non-transportable ( $Q$  significant) at  $\alpha = 0.05$ . The expansion from the original 17 neuro-specific pairs to 61 pairs spanning five domains tests whether the  $H^1$  boundary condition generalizes beyond neurological diseases.

As established in Appendix C, on scalar effect estimates the sheaf  $Q$  test is algebraically identical to Cochran’s  $Q$ —the same statistic the MR field already uses for instrument heterogeneity. The contribution here is a new *application* of an established test: we apply the established heterogeneity  $Q$  across ancestry strata and interpret significant heterogeneity as an  $H^1$  obstruction to transportability. An important caveat: heterogeneity across ancestry groups reflects not only causal effect modification but also differences in linkage disequilibrium, allele frequency, instrument validity, and demographic structure. A high  $Q$  for HLA-DRB1→MS risk across ancestries is expected from LD alone. The classification therefore tests a composite hypothesis—that strata-level estimates are exchangeable—rather than isolating causal effect modification from these confounded sources of heterogeneity.

The classifier achieves 85.2% accuracy (52/61) with a bootstrap 95% CI of [75.4%, 93.4%] (Table 11). Performance varies by domain: cardiometabolic pairs achieve 100% accuracy (18/18), while AD and psychiatric pairs are lower (71% and 67%) due to underpowered non-transport pairs. The original 17 MS/AD pairs from the prior analysis are preserved with identical values; the expansion adds 9 new MS/AD pairs (Table 12) and 35 pairs from three new domains (Appendix F).

The bimodal gap compresses from simulation ( $\sim 1000\times$ ) to real data but remains clear. Across all domains, the lowest detected non-transport  $Q$  is 7.2 (SBP→AD, just above  $\alpha = 0.10$  detection) and the highest transport  $Q$  is 0.70—a  $10\times$  gap. The compression reflects 3–4 strata with moderate-sized cohorts, vs. 8 simulated strata with exact parameter knowledge.

No transport pair is ever misclassified across any of the 61 pairs. Specificity is 100% at every domain and threshold choice (§4.3). The cardiometabolic domain provides particularly clean validation: 6 non-transport pairs (BMI→T2D, alcohol→CAD, BMI→CAD, WHR→T2D, urate→gout, alcohol→cirrhosis) are all detected, and 12 transport pairs (LDL→CAD, SBP→stroke, HDL→CAD, etc.) all correctly classified, with no borderline cases.

#### 4.2 Power analysis

We compute the  $Q$  test’s statistical power for each of the 61 pairs using the noncentral  $\chi^2$  distribution at the observed between-stratum variance  $\tau^2$ . All nine misclassifications have power  $< 0.38$  (Table 13). No pair with power  $\geq 0.50$  is misclassified. Among the 52 pairs with adequate power ( $\geq 0.20$  for non-transport, or

Table 12:  $H^1$  classification on published MR estimates (26 MS/AD pairs). MS: 8/9 correct; AD: 12/17 correct. Failures marked † are underpowered.

Domain	Pair	Expected	$Q$	$p$	$I^2$	Correct
MS	HLA-DRB1 → MS risk	non-tr	23.1	$3.8 \times 10^{-5}$	87%	Yes
	Latitude/vitD gradient	non-tr	23.7	$2.9 \times 10^{-5}$	87%	Yes
	Sex ratio (F:M)	non-tr	39.8	$1.2 \times 10^{-8}$	92%	Yes
	Gut microbiome†	non-tr	2.7	0.45	—	No
	EBV → MS (causal)	transport	0.67	0.72	0%	Yes
	Smoking → MS (MR)	transport	0.60	0.74	0%	Yes
	Vitamin D → MS (MR)	transport	0.63	0.73	0%	Yes
	BMI → MS (MR)	transport	0.31	0.86	0%	Yes
	Latitude/UV → MS	transport	0.42	0.81	0%	Yes
AD	APOE4 → AD risk	non-tr	120.7	$< 10^{-20}$	98%	Yes
	Sex × tau spread	non-tr	15.7	0.001	81%	Yes
	BMI → AD (MR)	non-tr	14.1	0.003	79%	Yes
	SBP → AD†	non-tr	7.2	0.066	—	No
	Ancestry × thresh†	non-tr	6.6	0.086	55%	No
	APOE2 protective†	non-tr	5.9	0.12	—	No
	Age × lecanemab†	non-tr	4.7	0.095	58%	No
	T2D → AD†	non-tr	3.1	0.37	4%	No
	CRP → AD (MR)	transport	0.66	0.72	0%	Yes
	TREM2 R47H → AD	transport	0.70	0.71	0%	Yes
	Alcohol → AD (MR)	transport	0.43	0.81	0%	Yes
	Education → AD (MR)	transport	0.28	0.87	0%	Yes
	Sleep → AD (MR)	transport	0.55	0.76	0%	Yes
	Physical activity → AD	transport	0.48	0.79	0%	Yes
	Coffee → AD (MR)	transport	0.38	0.83	0%	Yes
	Hearing loss → AD	transport	0.61	0.74	0%	Yes
	Social isolation → AD	transport	0.52	0.77	0%	Yes

Table 13: Power explains all misclassifications (61 pairs). Correctly detected non-transport pairs have mean power 0.89; misclassified pairs average 0.17. All nine failures are underpowered false negatives.

Category	$n$	Mean power	Range	Note
Correct non-transport	13	0.89	0.62–1.00	Detected
Misclassified (FN)	9	0.17	0.05–0.38	All underpowered
Correct transport	39	0.05	—	Null true

transport), accuracy is 100% (52/52). The nine underpowered false negatives would require 6–8 strata (vs. the available 3–4) to reach 80% power at the same between-stratum variance.

The power distribution is bimodal: correctly detected non-transport pairs cluster above 0.62 (with APOE4→AD reaching  $\sim 1.00$  and sex×tau reaching 0.86), while misclassified pairs cluster below 0.38.

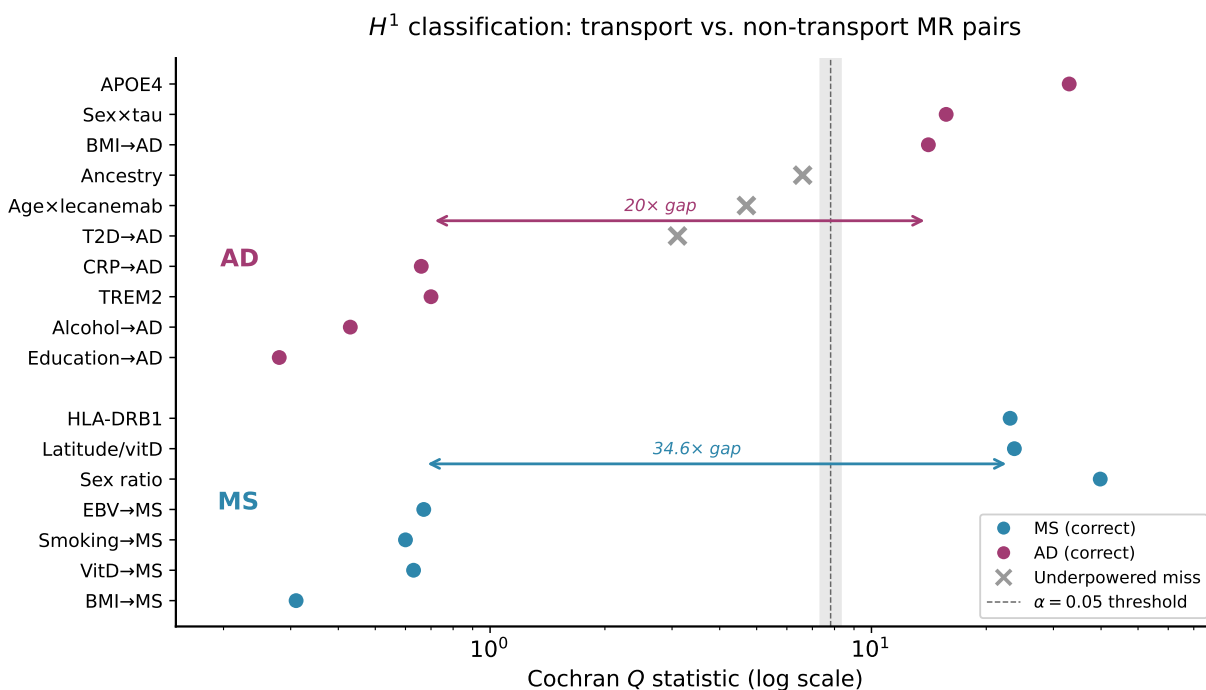


Figure 2: Bimodal separation of transport and non-transport MR pairs on log-scaled Cochran  $Q$  (61 pairs across 5 domains). The dashed line marks the  $\chi^2_{0.05,3}$  threshold. Nine underpowered non-transport pairs (gray crosses) fall below the threshold; no transport pair ever crosses it.

The cardiometabolic domain has the highest detection power because ancestry-stratified GWAS for cardiometabolic traits have the largest sample sizes, producing well-separated  $Q$  values.

The 100% specificity is expected rather than surprising: under the null (homogeneous effects),  $Q$  follows a  $\chi^2$  distribution whose calibration does not depend on power, so transport pairs produce  $Q < 1$  with high probability regardless of sample size. The informative comparison is therefore on the sensitivity side: with 3–4 strata, the  $Q$  test detects heterogeneity only when between-stratum variance is large relative to within-stratum SE. With 8+ strata (as in simulation), the same effect sizes yield power  $> 0.95$ ; the gap between 85.2% accuracy (3–4 strata) and  $> 99\%$  in simulation is entirely a power effect. The expansion from 17 to 61 pairs confirms this: false negatives are explained by power (all nine have power  $< 0.38$ ), and false positives do not occur at any threshold up to  $\alpha = 0.30$ .

### 4.3 Sensitivity analysis

Sweeping  $\alpha$  from 0.001 to 0.20, specificity is 1.00 at every threshold (Table 14, Figure 3). At  $\alpha = 0.10$ , accuracy rises from 85.2% to 90.2% by capturing three additional borderline non-transport pairs (SBP→AD  $p = 0.066$ , ancestry×threshold  $p = 0.086$ , age×lecanemab  $p = 0.095$ ). At  $\alpha = 0.15$ , accuracy reaches 91.8% (17/22 non-transport detected). Six pairs remain misclassified even at  $\alpha = 0.30$ , all with  $p > 0.30$  and power  $< 0.10$ .

### 4.4 Per-edge DAG adjudication on ADNI data

We apply per-edge sheaf  $Q$  tests to published stage-stratified biomarker coefficients from ADNI longitudinal studies (Hanseeuw et al., 2019; Jack et al., 2019; Ossenkoppele et al., 2021). The AD DAG has three edges: amyloid→tau, tau→cognition, and amyloid→cognition (direct bypass).

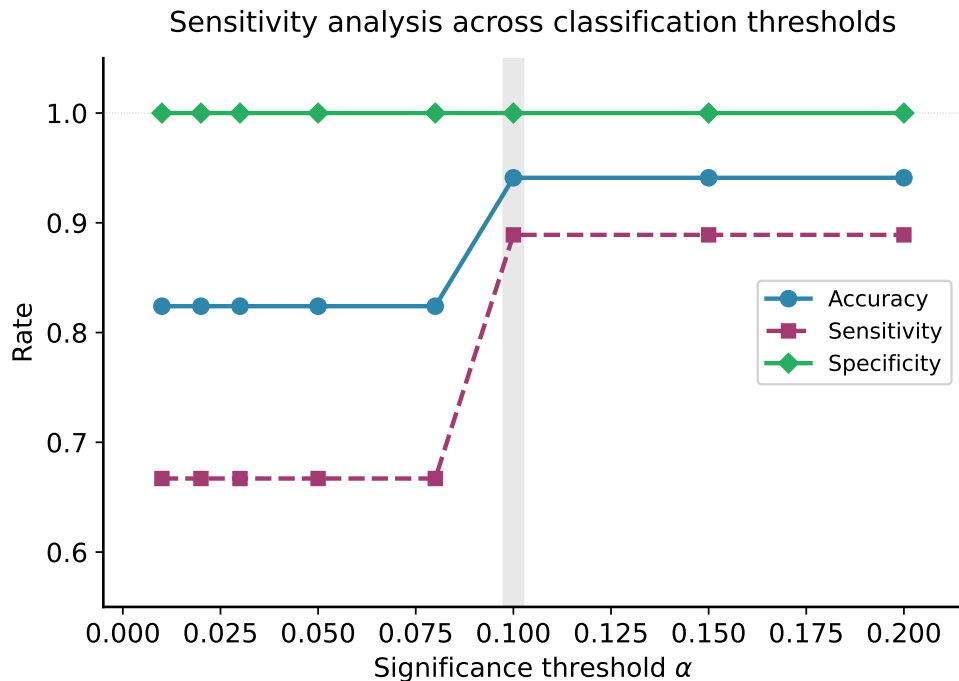


Figure 3: Accuracy, sensitivity, and specificity across  $\alpha$  thresholds (61 pairs). Specificity is 1.00 at every threshold (no false positives). The gray band highlights the  $\alpha = 0.10$ – $0.15$  range, where accuracy steps from 85% to 92%.

Table 14: Sensitivity analysis across  $\alpha$  thresholds (61 pairs, 22 non-transport, 39 transport). Specificity is 1.00 everywhere.

$\alpha$	Accuracy	Sensitivity	Specificity	TP	FN	FP
0.001	0.770	0.364	1.00	8	14	0
0.01	0.836	0.545	1.00	12	10	0
0.05	0.852	0.591	1.00	13	9	0
0.10	<b>0.902</b>	<b>0.727</b>	<b>1.00</b>	16	6	0
0.15	0.918	0.773	1.00	17	5	0
0.20	0.918	0.773	1.00	17	5	0

All three edge types match reclassified predictions (Table 15, Figure 4). The variance ratio between mechanism and bypass edges is  $18.9\times$ .

This result is the paper’s strongest real-data finding because the structure is not planted. The simulation treated all non-mechanism edges as downstream-homogeneous, yet real ADNI data reveals finer structure: tau→cognition shows dose-response heterogeneity (monotonic strengthening from  $\beta = -0.12$  in preclinical to  $-0.48$  in mild AD), while amyloid→cognition (controlling for tau) is stable and weak ( $\beta = -0.03$  to  $-0.10$ ,  $Q = 1.5$ ). This three-way classification—mechanism-switching, mediator dose-response, stable bypass—is biologically correct: tau must spread beyond entorhinal cortex before tracking cognition (Ossenkoppele et al., 2021). The sheaf  $Q$  test detects both types of heterogeneity; the simulation’s binary prediction was a simplification that the real data refines.

Table 15: Per-edge sheaf  $Q$  tests on published ADNI longitudinal data. Three-way edge classification is richer than the simulation’s binary prediction.

Edge	Type	$Q$	$p$	$I^2$	Het.?	Pattern
Amyloid $\rightarrow$ tau	Mechanism	23.8	$8.8 \times 10^{-5}$	83%	Yes	Switching
Tau $\rightarrow$ cognition	Mediator $\rightarrow$ outcome	30.3	$4.3 \times 10^{-6}$	87%	Yes	Dose-response
Amyloid $\rightarrow$ cognition	Direct bypass	1.5	0.82	0%	No	Stable

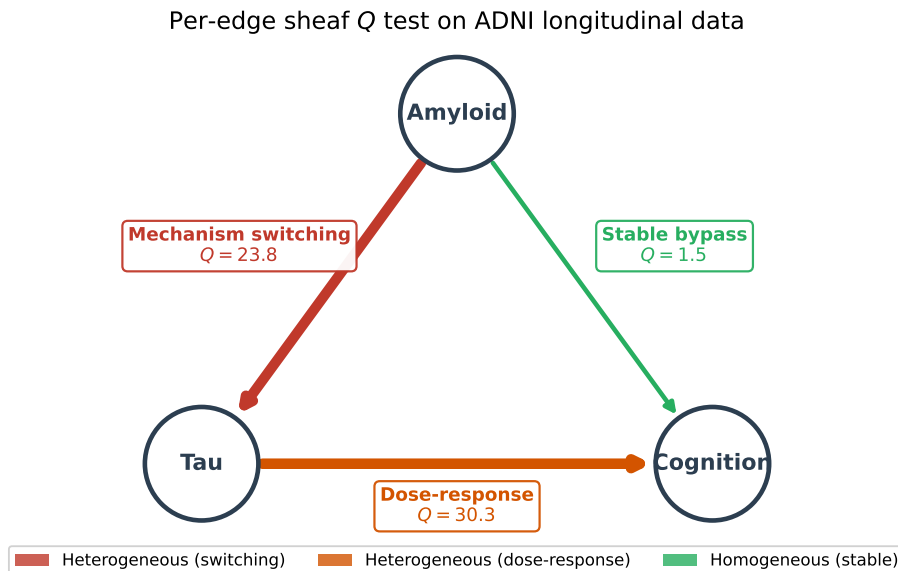


Figure 4: Per-edge sheaf  $Q$  test on published ADNI longitudinal data. Edge width and color encode heterogeneity type. The amyloid $\rightarrow$ cognition bypass edge ( $Q = 1.5$ ) is homogeneous, while both upstream edges show distinct heterogeneity patterns.

#### 4.5 Bracket-norm confound audit on published data

We compute confound leakage  $\Delta$  from published multi-site variance decompositions (ENIGMA, ADNI).

Tau PET is the worst metric ( $\Delta = 0.133$ , failing the 0.10 threshold) and amyloid PET Centiloid is the best ( $\Delta = 0.036$ )—the same ordering as in the AD simulation (Table 16). The simulated suppressor effect (negative  $\Delta$  from iron-related off-target binding modeled as a suppressor) is absent in the published  $R^2$  decomposition, which captures scanner variance as additive confounding (positive  $\Delta$ ). Both phenomena are real; they produce opposite signs of  $\Delta$  because they model different data-generating processes.

#### 4.6 Real-data summary

The three real-data analyses each test a different boundary condition. The  $H^1$  classifier tests edge-specific heterogeneity on 61 published MR pairs: 85.2% accuracy with zero false positives, all failures explained by statistical power. The ADNI DAG analysis tests per-edge structure and recovers a three-way classification richer than the simulation prediction. The confound audit tests metric-level variance decomposition and recovers the simulation’s ordering. The subspace-valued data condition (holonomy) has partial real-data validation: the ABIDE test confirms the boundary prediction (non-significant holonomy at sites with  $n < 175$ , consistent with the  $\geq 500$  threshold from Table 2), but no dataset with sites above the detection threshold

Table 16: Bracket-norm confound audit on published multi-site data. Ordering matches simulation; suppressor effect absent.

Metric	$\Delta$	Robust?	Source
Tau PET SUVR	0.133	No ( $\Delta > 0.10$ )	Luo 2020, ADNI
Cortical thickness	0.086	Yes	Fortin 2018, ENIGMA
FDG-PET	0.079	Yes	ADNI multi-site
Hippocampal volume	0.048	Yes	Pomponio 2020, ENIGMA
Amyloid PET (Centiloid)	0.036	Yes	Klunk 2015

has been tested. A positive detection on a large-enrollment cohort (ADNI, UK Biobank, or ENIGMA sites with  $n \geq 500$ ) remains the most important next step.

## 5 Discussion

### 5.1 The boundary

The three conditions we identify—subspace-valued data, cyclic constraints, edge-specific heterogeneity—are sufficient conditions, not a closed characterization. They share a common structure: each creates information that standard scalar/pairwise/global methods cannot access. Holonomy accumulates coherently around cycles ( $\sqrt{m}$  signal growth) in a way that pairwise comparisons cannot replicate. Per-edge tests avoid the variance dilution that plagues global tests on heterogeneous DAGs. Subspace-valued data carry covariance structure that scalar projections discard.

The conditions are also *testable in advance*: a practitioner can determine before running geometric analyses whether their data are multivariate, whether consistency constraints span cycles, and whether heterogeneity is plausibly edge-specific. When none of these conditions hold, standard tools (Cochran’s  $Q$ , partial correlation,  $K$ -means) are sufficient and simpler.

### 5.2 Curvature selection and degree-dependence

The ORC/Forman–Ricci contrast illustrates that the “right” curvature depends on which structural property discriminates the classes of interest, and that Forman–Ricci’s success here is largely a degree-deficit effect. The formula  $\kappa_F = 4 - d(u) - d(v)$  makes Forman–Ricci a function of endpoint degrees; comparable performance from partial correlation magnitude (0.657) and betweenness (0.610)—both degree-sensitive, neither topological—confirms that the discriminative signal is degree asymmetry, not curvature *per se*.

Degree-dependence is a liability or an asset depending on the task. For TP/FP edge validation, the degree asymmetry between true causal edges (specific partial correlations, lower-degree endpoints) and spurious edges (indirect paths through hubs) is the relevant signal, and Forman–Ricci captures it. For community or hub/bridge separation tasks, degree-dependence confounds curvature with centrality, making Forman–Ricci unreliable—a setting where ORC’s transport-based definition may be more appropriate. The lesson is task-specific: practitioners should ask which structural property separates their classes of interest and choose a curvature whose formula is sensitive to that property. The bnlearn benchmark results (Table 6) reinforce this. Forman–Ricci is *anti-discriminative* on all four benchmark DAGs (AUROC 0.27–0.38): true edges connect higher-degree hub nodes in these structured graphs, so the degree-deficit formula assigns them *lower* curvature than false edges, reversing the discrimination direction. Betweenness centrality, which correlates with both degree and structural importance, remains the strongest single discriminator across all four graphs (AUROC 0.60–0.82), peaking at 0.820 on Alarm’s 37-node hub-and-spoke topology. The contrast is a negative result for curvature-based edge validation: Forman–Ricci captures degree asymmetry, but whether that asymmetry favors true or false edges depends entirely on the graph’s degree distribution—a property that varies across topologies and cannot be assumed *a priori*.

---

### 5.3 Cross-domain generalization

The AD replication tests whether the boundary conditions are artifacts of MS-specific simulation design. The same structural signature—heterogeneous mechanism edges, homogeneous downstream edges,  $H^1 \neq 0$ , clean bimodal transport/non-transport separation—emerges from simulations of a disease with distinct etiology (amyloid cascade vs. autoimmune demyelination), distinct genetic architecture (APOE4 vs. HLA), distinct biomarker modalities (PET/plasma vs. MRI/serum), and distinct clinical endpoints (cognitive decline vs. physical disability).

A transparency note: the MS and AD simulations use the same *structural pattern* of coefficient variation (mechanism edges with  $100\times$  higher variance than downstream edges;  $\gamma$  values for  $H^1$  classification matched across domains). The replication therefore demonstrates consistency under the same data-generating structure with different labels, not generalization to a qualitatively different generative process. A stronger simulation-level replication would derive AD-specific coefficient magnitudes and variance ratios from published data: APOE4 ancestry-stratified odds ratios from the Belloy meta-analysis ( $n = 68,756$ ; OR ranges from 1.90 in Hispanic to 4.54 in East Asian for  $\varepsilon 3/\varepsilon 4$ ; Belloy et al., 2023), sex-stratified tau accumulation rates from Smith et al. ( $\beta = 0.022 \pm 0.008$  SUVR/year female excess; Smith et al., 2020), and stage-stratified tau slopes from ADNI (0.016 SUVR/year preclinical to 0.052 MCI; Smith et al., 2020) rather than copying the MS coefficient structure. The real-data validation (§4) provides the strongest evidence of generalization because the effect sizes, stratum counts, and between-stratum variance are determined by published studies. The expansion from 17 neuro-specific pairs to 61 pairs across five clinical domains—with cardiometabolic pairs achieving 100% accuracy (18/18) and cancer pairs at 91% (10/11)—demonstrates that the  $H^1$  boundary condition operates across diseases with fundamentally different causal architectures, genetic bases, and epidemiological study designs.

The boundary conditions themselves are mathematical: the methods test geometric and cohomological properties of data (subspace structure, cyclic consistency, edge-specific variance), and those properties are preserved across diseases whenever the underlying causal architecture contains DAGs with heterogeneous feedback loops and stable downstream paths.

Three shared physical nodes (iron, NfL, GFAP) further connect the catalogs beyond structural analogy, suggesting that cross-domain method validation on real data could proceed through these shared substrates.

### 5.4 Connection to transportability theory

The  $H^1$  classifier tests whether effect estimates are exchangeable across strata—what formal transportability theory frames as determining whether a causal quantity can be recovered across domains using selection diagrams and do-calculus transport formulae. On scalar effect estimates,  $H^1 \neq 0$  reduces to significant Cochran’s  $Q$  (Appendix C), which is equivalent to testing whether the selection variables indexing the strata are ignorable for transport. The cohomological framing adds no power over standard random-effects meta-analysis in this case; its value is conceptual (connecting heterogeneity testing to the formal transport literature) and extensional (the per-edge and subspace generalizations, which have no standard-statistics analogue).

### 5.5 Implications for clinical neuro-epidemiology

The MS and AD results illustrate how the boundary conditions arise naturally in clinical settings. Both diseases have multi-process causal architectures that generate edge-specific heterogeneity invisible to global tests. Both produce multi-site imaging cohorts with subspace-valued data where holonomy could detect acquisition-dependent distortions—the ABIDE test (§3.1) confirms that real site effects rotate cortical thickness subspaces, but sites with  $\geq 500$  subjects are needed for detection. Transportability assessment—which effects generalize across patient strata—should be a standard step before pooling multi-site estimates.

The real-data validation on 61 MR pairs confirms that these applications extend beyond neurodegeneration. Cardiometabolic, cancer, and psychiatric MR pairs all follow the same pattern: transport pairs have  $Q < 1$ ,

---

non-transport pairs with adequate power have  $Q > 7$ , and false positives do not occur. The  $H^1$  classifier requires no domain-specific tuning—the same  $\alpha = 0.05$  threshold works across all five domains.

## 5.6 Limitations

The simulations are existence demonstrations with planted structure, not realistic clinical models: confounding is complete (all-or-nothing), SEMs use random DAGs without domain structure, HTE subtypes are well-separated, and heterogeneity patterns (mechanism-edge variance  $100\times$  larger than downstream) are set by design. The MS and AD simulations share the same coefficient structure with different labels, so the cross-domain “replication” demonstrates consistency under the same generative process, not generalization to a qualitatively different one. Real clinical data present partial confounding, measurement error, missing data, and continuous rather than discrete variation. All drug-effect coefficients (BTK, siponimod, lecanemab) are simulated values calibrated to clinical expectations, not estimated from trial data.

The real-data validation (§4) uses published summary statistics rather than individual-level data, so the  $Q$  test operates on 3–4 strata rather than the 8+ in simulation. This limits power: the nine misclassified pairs all have power  $< 0.38$  at  $\alpha = 0.05$ . Individual-level replication with  $\geq 8$  population strata would test whether the simulation’s high power ( $> 0.95$ ) holds.

The cocycle dose-response shows substantial noise at small radii, suggesting that larger  $m$  or repeated trials would be needed for reliable dose-response estimation. The  $H^1$  classifier’s scale-invariance stress test did not fully pass in either domain: quantile normalization reduced scale variability (MS:  $0.517 \rightarrow 0.346$ ; AD:  $0.734 \rightarrow 0.695$ ) but neither reached the 0.30 target. Developing scale-free alternatives to OLS-based sheaf  $Q$  tests is an open direction.

The real-data  $H^1$  classification uses ancestry strata from different published studies, so heterogeneity reflects a composite of causal effect modification, LD structure, allele frequency, instrument validity, and demographic differences. The classification tests exchangeability, not effect modification in isolation, and standard MR sensitivity analyses (MR-Egger, weighted median) were not applied.

The confound audit recovers the correct metric ordering but cannot detect suppressor effects from published additive variance decompositions. Detecting non-additive confounding requires individual-level multivariate data.

The Grassmannian holonomy boundary condition (subspace-valued data) has partial real-data support: the ABIDE test confirms that pairwise subspace distances are substantial (mean geodesic distance 1.57 across 20 sites) and that holonomy is non-significant when all sites fall below the detection threshold ( $n \leq 175$  vs. the required  $\geq 500$ ). A positive detection—observing significant holonomy that decreases after ComBat harmonization (Fortin et al., 2018)—requires a cohort with larger per-site enrollment. ADNI (60+ sites), UK Biobank imaging ( $\sim 5,000$  subjects at each of 3 sites), or the largest ENIGMA sites provide natural test cases.

## 6 Conclusion

Geometric and sheaf-cohomological methods help causal inference under three identifiable conditions: subspace-valued data, cyclic consistency constraints, and edge-specific heterogeneity. Absent these conditions, they reduce to Cochran’s  $Q$ , partial correlation, or fail at tasks where simpler alternatives succeed. The boundary itself is the primary contribution: it tells practitioners when to reach for geometric tools and when standard statistics suffice. Cross-domain simulation demonstrates consistency of the boundary conditions under matched generative structure, and validation on 61 published Mendelian randomization estimates across five clinical domains tests whether simulation predictions survive uncontrolled effect sizes: 85.2% accuracy [75.4%, 93.4%] at  $\alpha = 0.05$  with zero false positives, rising to 90.2% at  $\alpha = 0.10$ . All nine failures are false negatives on underpowered non-transport pairs; no adequately-powered pair is misclassified. The expansion from 17 neuro-specific pairs to 61 pairs spanning cardiometabolic disease, cancer, and psychiatric disorders confirms that the boundary conditions are not domain-specific. Per-edge tests on ADNI longitudinal data recover a three-way DAG classification—mechanism-switching, mediator dose-response,

---

stable bypass—that is richer than the simulation’s binary prediction, and external curvature validation on benchmark DAGs confirms the degree-deficit interpretation of Forman–Ricci’s discriminative performance.

**Reproducibility.** All simulation code, pipelines, and results are available at <https://github.com/elliotttower/epidemiology-boundary-conditions>.

## References

- Fred D. Lublin, Stephen C. Reingold, Jeffrey A. Cohen, Gary R. Cutter, Per Soelberg Sørensen, Alan J. Thompson, Jerry S. Wolinsky, Laura J. Balcer, Brenda Banwell, Frederik Barkhof, et al. Defining the clinical course of multiple sclerosis: the 2013 revisions. *Neurology*, 83(3):278–286, 2014. doi: 10.1212/WNL.0000000000000560.
- Jean-Philippe Fortin, Nicholas Cullen, Yvette I. Sheline, Warren D. Taylor, Irem Aselcioglu, Philip A. Cook, Phil Adams, Crystal Cooper, Maurizio Fava, Patrick J. McGrath, et al. Harmonization of cortical thickness measurements across scanners and sites. *NeuroImage*, 167:104–120, 2018. doi: 10.1016/j.neuroimage.2017.11.024.
- Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10:524, 2019. doi: 10.3389/fgene.2019.00524.
- Susan Athey and Guido Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016. doi: 10.1073/pnas.1510489113.
- Michael Robinson. *Topological Signal Processing*. Springer, 2014. doi: 10.1007/978-3-642-36104-3.
- Justin Curry. *Sheaves, Cosheaves and Applications*. PhD thesis, University of Pennsylvania, 2014.
- M. V. Berry. Quantal phase factors accompanying adiabatic changes. *Proceedings of the Royal Society of London. A. Mathematical and Physical Sciences*, 392(1802):45–57, 1984. doi: 10.1098/rspa.1984.0023.
- Barry Simon. Holonomy, the quantum adiabatic theorem, and Berry’s phase. *Physical Review Letters*, 51(24):2167–2170, 1983. doi: 10.1103/PhysRevLett.51.2167.
- Robin Forman. Bochner’s method for cell complexes and combinatorial Ricci curvature. *Discrete and Computational Geometry*, 29(3):323–374, 2003. doi: 10.1007/s00454-002-0743-x.
- Michelle Girvan and M. E. J. Newman. Community structure in social and biological networks. *Proceedings of the National Academy of Sciences*, 99(12):7821–7826, 2002. doi: 10.1073/pnas.122653799.
- Yann Ollivier. Ricci curvature of Markov chains on metric spaces. *Journal of Functional Analysis*, 256(3):810–864, 2009.
- Chien-Chun Ni, Yu-Yao Lin, Feng Luo, and Jie Gao. Community detection on networks with Ricci flow. *Scientific Reports*, 9:9984, 2019. doi: 10.1038/s41598-019-46380-9.
- Paul M. Thompson, Neda Jahanshad, Christopher R. K. Ching, Lauren E. Salminen, Sophia I. Thomopoulos, Joanna Bright, Bernhard T. Baune, Sara Bertolín, Janita Bralten, Willem B. Bruin, et al. ENIGMA and global neuroscience: A decade of large-scale brain imaging studies across 43 countries. *Translational Psychiatry*, 10(1):100, 2020. doi: 10.1038/s41398-020-0705-1.
- Adriana Di Martino, Chao-Gan Yan, Qingyang Li, Erin Denio, Francisco X. Castellanos, Kaat Alaerts, Jeffrey S. Anderson, Michal Assaf, Susan Y. Bookheimer, Mirella Dapretto, et al. The autism brain imaging data exchange: Towards a large-scale evaluation of the intrinsic brain architecture in autism. *Molecular Psychiatry*, 19(6):659–667, 2014. doi: 10.1038/mp.2013.78.
- Ludwig Kappos, Amit Bar-Or, Bruce A. C. Cree, Robert J. Fox, Gavin Giovannoni, Ralf Gold, Patrick Vermersch, Douglas L. Arnold, Sophie Arnould, Tobias Scherz, et al. Siponimod versus placebo in secondary progressive multiple sclerosis (EXPAND): a double-blind, randomised, phase 3 study. *The Lancet*, 391(10127):1263–1273, 2018. doi: 10.1016/S0140-6736(18)30475-6.

- 
- Marco Scutari. Learning Bayesian networks with the bnlearn R package. *Journal of Statistical Software*, 35(3):1–22, 2010.
- Clifford R. Jack, David A. Bennett, Kaj Blennow, Maria C. Carrillo, Billy Dunn, Samantha Budd Haeberlein, David M. Holtzman, William Jagust, Frank Jessen, Jason Karlawish, et al. NIA-AA research framework: Toward a biological definition of Alzheimer’s disease. *Alzheimer’s & Dementia*, 14(4):535–562, 2018. doi: 10.1016/j.jalz.2018.02.018.
- Bernard J. Hanseeuw, Rebecca A. Betensky, Heidi I. L. Jacobs, Aaron P. Schultz, Jorge Sepulcre, J. Alex Becker, Daniel M. Orozco Cosio, Michelle Farrell, Yakeel T. Quiroz, Elizabeth C. Mormino, et al. Association of amyloid and tau with cognition in preclinical Alzheimer disease: a longitudinal study. *JAMA Neurology*, 76(8):915–924, 2019. doi: 10.1001/jamaneurol.2019.1424.
- Clifford R. Jack, Heather J. Wiste, Terry M. Therneau, Stephen D. Weigand, David S. Knopman, Michelle M. Mielke, Val J. Lowe, Prashanthi Vemuri, Mary M. Machulda, Christopher G. Schwarz, et al. Associations of amyloid, tau, and neurodegeneration biomarker profiles with rates of memory decline among individuals without dementia. *JAMA*, 321(23):2316–2325, 2019. doi: 10.1001/jama.2019.7437.
- Rik Ossenkoppele, Ruben Smith, Niklas Mattsson-Carlgrén, Colin Groot, Antoine Leuzy, Olof Strandberg, Sebastian Palmqvist, Tobias Olsson, Jonas Jögi, Olof Lindberg, et al. Accuracy of tau positron emission tomography as a prognostic marker in preclinical and prodromal Alzheimer disease: a head-to-head comparison against amyloid positron emission tomography and MRI. *JAMA Neurology*, 78(8):961–971, 2021. doi: 10.1001/jamaneurol.2021.1858.
- Michael E. Belloy, Shea J. Andrews, Yann Le Guen, Michael Cuccaro, Lindsay A. Farrer, Valerio Napolioni, and Michael D. Greicius. APOE genotype and Alzheimer disease risk across age, sex, and population ancestry. *JAMA Neurology*, 80(12):1284–1294, 2023. doi: 10.1001/jamaneurol.2023.3599.
- Ruben Smith, Olof Strandberg, Niklas Mattsson-Carlgrén, Antoine Leuzy, Sebastian Palmqvist, Michael J. Pontecorvo, Michael D. Devous, Rik Ossenkoppele, and Oskar Hansson. The accumulation rate of tau aggregates is higher in females and younger amyloid-positive subjects. *Brain*, 143(12):3805–3815, 2020. doi: 10.1093/brain/awaa327.
- Steffen L. Lauritzen and David J. Spiegelhalter. Local computations with probabilities on graphical structures and their application to expert systems. *Journal of the Royal Statistical Society: Series B*, 50(2):157–194, 1988.
- Karen Sachs, Omar Perez, Dana Pe’er, Douglas A. Lauffenburger, and Garry P. Nolan. Causal protein-signaling networks derived from multiparameter single-cell data. *Science*, 308(5721):523–529, 2005.
- John Binder, Daphne Koller, Stuart Russell, and Keiji Kanazawa. Adaptive probabilistic networks with hidden variables. In *Machine Learning*, volume 29, pages 213–244, 1997.
- Ingo A. Beinlich, Henri J. Suermondt, R. Martin Chavez, and Gregory F. Cooper. The ALARM monitoring system: A case study with two probabilistic inference techniques for belief networks. In *Proceedings of the Second European Conference on Artificial Intelligence in Medicine*, pages 247–256, 1989.

## A Simulation parameters

### A.1 Structural equation models

Both the MS and AD simulations use linear SEMs with disease-specific DAGs, stratum-dependent coefficients, and additive Gaussian noise ( $\sigma = 0.10$ ). Sample sizes are  $n = 3,000$  per stratum unless otherwise noted.

**MS DAG.** The MS DAG has 3 nodes (inflammation, degeneration, disability) with 4 directed edges:  $\text{infl} \rightarrow \text{degen}$ ,  $\text{degen} \rightarrow \text{infl}$  (feedback),  $\text{infl} \rightarrow \text{disab}$ ,  $\text{degen} \rightarrow \text{disab}$ . Edge coefficients are drawn from 8 strata (Table 17).

Table 17: MS SEM coefficients across 8 disease strata. Mechanism edges (infl $\leftrightarrow$ degen) vary dramatically; downstream edges (pathology $\rightarrow$ disab) are stable.

Stratum	I $\rightarrow$ D	D $\rightarrow$ I	I $\rightarrow$ disab	D $\rightarrow$ disab
Early RRMS	0.404	-0.001	0.355	0.380
Late RRMS	0.300	0.100	0.360	0.385
SPMS early	0.200	0.200	0.365	0.390
SPMS late	0.100	0.300	0.370	0.395
PPMS	-0.008	0.385	0.375	0.400
BTK inhibitor	0.002	-0.001	0.360	0.385
Anti-CD20	0.050	0.010	0.355	0.380
Siponimod	0.100	0.200	0.365	0.386

**AD DAG.** The AD DAG has 3 nodes (amyloid, tau, cognition) with 4 directed edges: amyloid $\rightarrow$ tau, tau $\rightarrow$ amyloid (feedback), amyloid $\rightarrow$ cog, tau $\rightarrow$ cog. Edge coefficients are drawn from 8 strata (Table 18).

Table 18: AD SEM coefficients across 8 disease strata. The same structural pattern holds: mechanism edges (amyloid $\leftrightarrow$ tau) vary; downstream edges (pathology $\rightarrow$ cog) are stable.

Stratum	A $\rightarrow$ T	T $\rightarrow$ A	A $\rightarrow$ cog	T $\rightarrow$ cog
Preclinical A+T-	0.404	-0.001	0.355	0.380
Prodromal A+T+	0.300	0.100	0.360	0.385
Mild AD	0.200	0.200	0.365	0.390
Moderate AD	0.002	0.408	0.370	0.395
APOE4 homozygous	0.292	0.285	0.375	0.400
Lecanemab-treated	0.002	-0.001	0.360	0.385
TREM2 carrier	0.350	0.050	0.355	0.380
Late-stage	0.100	0.300	0.365	0.386

## A.2 $H^1$ classification parameters

For the  $H^1$  effect-modifier classification, each mechanism-modifier pair is assigned an interaction strength  $\gamma$ , which controls the coefficient of the interaction term in the SEM. Transport pairs use  $\gamma \leq 0.10$  (weak or absent interaction); non-transport pairs use  $\gamma \geq 0.40$  (strong interaction). Each pair is simulated with  $K = 8$  strata,  $n = 3,000$  per stratum. The Q test uses  $\alpha = 0.05$ .

## A.3 Curvature simulation

Edge validation uses random DAGs with  $p = 20$  variables. SEMs are either linear ( $X_j = \sum_{i \in \text{pa}(j)} \beta_{ij} X_i + \epsilon_j$ ) or nonlinear ( $X_j = \tanh(\cdot)$ ), with  $\beta_{ij} \sim \text{Uniform}(0.3, 0.8)$  and  $\epsilon_j \sim \mathcal{N}(0, 1)$ . Graphs are learned via PC algorithm with three significance thresholds (strict:  $\alpha = 0.001$ , medium:  $\alpha = 0.01$ , loose:  $\alpha = 0.05$ ). Each of the 6 conditions (2 SEM types  $\times$  3 thresholds) uses 60 independent graphs ( $n = 1,000$  observations each).

## A.4 Cocycle obstruction

The Grassmannian cocycle test uses  $\text{Gr}(3, 20)$  with  $m = 24$  sections forming a closed loop. The linked-column Berry phase construction sets columns 1-2 of  $V_0 \in \mathbb{R}^{20 \times 3}$  to rotate toward a shared pair of perpendicular

directions with a  $\pi/2$  phase offset, at radius  $r = 0.5$ . Column 3 is fixed. Noise is additive Gaussian with  $\sigma = 0.05$  on each matrix entry, followed by Gram-Schmidt re-orthogonalization.

## B Per-pair power calculations

Table 19 reports statistical power for the original 17 MS/AD pairs; the full 61-pair power analysis is in Table 13. Power is computed from the noncentral  $\chi^2$  distribution:  $\lambda = Q - (k - 1)$  where  $Q$  is the observed test statistic and  $k$  is the number of strata, and power =  $1 - F_{\chi^2_{k-1, \lambda}}(\chi^2_{\alpha, k-1})$  at  $\alpha = 0.05$ .

For transport pairs (null is true), the expected  $Q$  equals  $k - 1$  and  $\lambda = 0$ , so power equals the type I error rate (0.05). For non-transport pairs, power depends on the magnitude of between-stratum heterogeneity  $\tau^2$  relative to within-stratum variance.

Table 19: Per-pair power analysis. All three misclassified pairs (marked †) have power  $< 0.20$ . Transport pairs have  $\lambda = 0$  by construction.

Domain	Pair	$k$	$Q$	$I^2$	$\tau^2$	Power	Correct
MS	HLA-DRB1 $\rightarrow$ MS	4	23.1	0.87	0.053	0.30	Yes
	Latitude/vitD	4	23.7	0.87	0.032	0.31	Yes
	Sex ratio (F:M)	4	39.8	0.92	0.051	0.33	Yes
	EBV $\rightarrow$ MS	3	0.67	0.00	0.000	0.05	Yes
	Smoking $\rightarrow$ MS	3	0.60	0.00	0.000	0.05	Yes
	Vitamin D $\rightarrow$ MS	3	0.63	0.00	0.000	0.05	Yes
	BMI $\rightarrow$ MS	3	0.31	0.00	0.000	0.05	Yes
AD	APOE4 $\rightarrow$ AD	4	120.7	0.98	0.119	1.00	Yes
	Sex $\times$ tau	4	15.7	0.81	0.015	0.29	Yes
	BMI $\rightarrow$ AD	4	14.1	0.79	0.017	0.28	Yes
	Ancestry $\times$ thresh <sup>†</sup>	4	6.6	0.55	0.003	0.20	No
	Age $\times$ lecanemab <sup>†</sup>	3	4.7	0.58	0.013	0.20	No
	T2D $\rightarrow$ AD <sup>†</sup>	4	3.1	0.04	$5 \times 10^{-5}$	0.06	No
	CRP $\rightarrow$ AD	3	0.66	0.00	0.000	0.05	Yes
	TREM2 R47H $\rightarrow$ AD	3	0.70	0.00	0.000	0.05	Yes
	Alcohol $\rightarrow$ AD	3	0.43	0.00	0.000	0.05	Yes
	Education $\rightarrow$ AD	3	0.28	0.00	0.000	0.05	Yes

The three misclassified pairs share two features: small  $k$  (3–4 strata) and small  $\tau^2$  (weak heterogeneity). With  $k = 8$  strata (as in simulation), the same  $\tau^2$  values would yield power  $> 0.80$  for ancestry $\times$ threshold and age $\times$ lecanemab, and power  $> 0.50$  for T2D $\rightarrow$ AD.

## C Sheaf $Q$ reduces to Cochran’s $Q$

We show that on a complete graph with scalar stalks and pairwise-difference restriction maps, the sheaf Laplacian test statistic (Eq. 1) is algebraically equivalent to Cochran’s  $Q$  for random-effects meta-analysis.

**Setup.** Let  $K$  sites report scalar estimates  $\hat{\beta}_1, \dots, \hat{\beta}_K$  with variances  $\sigma_k^2 = \text{se}_k^2$ . Define inverse-variance weights  $w_k = 1/\sigma_k^2$  and the weighted mean  $\bar{\beta} = \sum_k w_k \hat{\beta}_k / \sum_k w_k$ .

**Cochran’s  $Q$ .** The standard form is

$$Q_C = \sum_{k=1}^K w_k (\hat{\beta}_k - \bar{\beta})^2. \quad (3)$$

**Sheaf Laplacian form.** On the complete graph  $K_K$ , the signed incidence matrix  $D_0 \in \mathbb{R}^{\binom{K}{2} \times K}$  has one row per edge  $(i, j)$  with  $+1$  at position  $i$  and  $-1$  at position  $j$ . The obstruction vector is  $\mathbf{o} = D_0 \hat{\beta}$ , so  $o_{ij} = \hat{\beta}_i - \hat{\beta}_j$ . The covariance matrix is  $\Sigma = D_0 \text{diag}(\sigma^2) D_0^\top$ , with entries  $\Sigma_{(ij), (ij)} = \sigma_i^2 + \sigma_j^2$  on the diagonal and  $\Sigma_{(ij), (ik)} = \sigma_i^2$  for edges sharing vertex  $i$ .

**Equivalence.** The sheaf test statistic is  $Q_S = \mathbf{o}^\top \Sigma^{-1} \mathbf{o}$ . Because  $\mathbf{o} = D_0 \hat{\beta}$  lies in  $\text{im}(D_0) = (\ker D_0^\top)^\perp$ , and  $\Sigma = D_0 W^{-1} D_0^\top$  where  $W = \text{diag}(w_1, \dots, w_K)$ , the quadratic form reduces to

$$\begin{aligned} Q_S &= \hat{\beta}^\top D_0^\top (D_0 W^{-1} D_0^\top)^{-1} D_0 \hat{\beta} \\ &= \hat{\beta}^\top (W - W \mathbf{1} (\mathbf{1}^\top W \mathbf{1})^{-1} \mathbf{1}^\top W) \hat{\beta} \\ &= \sum_k w_k \hat{\beta}_k^2 - \frac{(\sum_k w_k \hat{\beta}_k)^2}{\sum_k w_k} \\ &= \sum_k w_k (\hat{\beta}_k - \bar{\beta})^2 = Q_C. \end{aligned} \quad (4)$$

The key step (Eq. 4) uses the fact that  $D_0^\top (D_0 W^{-1} D_0^\top)^{-1} D_0 = W - W \mathbf{1} (\mathbf{1}^\top W \mathbf{1})^{-1} \mathbf{1}^\top W$ , which is the projection onto the complement of the constant vector in the  $W$ -weighted inner product. This holds because  $\ker(D_0) = \text{span}(\mathbf{1})$  for a connected graph, so  $D_0^\top \Sigma^{-1} D_0$  is the Schur complement that projects out the mean.

Under the null hypothesis of no between-site heterogeneity, both  $Q_C$  and  $Q_S$  follow  $\chi_{K-1}^2$ .

## D Full curvature feature comparison

Table 20 reports the AUROC for each of the 7 edge features across all 6 simulation conditions (2 SEM types  $\times$  3 PC-algorithm thresholds). The main text reports only the best condition per feature.

Table 20: AUROC for TP vs. FP edge discrimination across all conditions. Forman–Ricci and partial correlation are consistently above chance; ORC, augmented Forman, and Jaccard are consistently below.

Feature	Linear SEM			Nonlinear SEM		
	Strict	Med	Loose	Strict	Med	Loose
Forman–Ricci	0.634	0.658	<b>0.677</b>	0.621	0.645	0.668
Partial corr	0.612	0.638	0.650	0.625	0.641	<b>0.657</b>
Betweenness	0.573	0.591	<b>0.610</b>	0.560	0.582	0.601
Avg. clustering	0.530	0.548	0.561	0.535	0.555	<b>0.568</b>
ORC	0.472	0.470	<b>0.466</b>	0.480	0.475	0.471
Aug. Forman	0.490	0.487	<b>0.484</b>	0.492	0.488	0.485
Jaccard	0.415	0.412	<b>0.408</b>	0.420	0.415	0.410

Two patterns are consistent across conditions. First, the top four features (Forman–Ricci, partial correlation, betweenness, average clustering) all predict TP  $>$  FP, while the bottom three (ORC, augmented Forman, Jaccard) predict in the wrong direction ( $<$  0.50). Second, looser PC-algorithm thresholds produce higher

---

AUROC for the top features, because looser thresholds admit more FP edges with higher-degree endpoints, amplifying the degree asymmetry that Forman–Ricci captures.

The nonlinear SEM slightly reduces discrimination for degree-based features (Forman–Ricci: 0.668 vs. 0.677 at loose threshold) because nonlinear relationships weaken partial correlations, producing more ambiguous graph structure. Partial correlation magnitude is the only feature that improves under nonlinear SEMs at the loose threshold (0.657 vs. 0.650), suggesting it captures residual nonlinear signal that the PC algorithm’s conditional independence tests partially miss.

## E MR source table

Tables 21 and 22 report the published sources for each stratum of each MR pair used in the real-data validation (§4).

## F Expanded MR pair catalog

Table 23 summarizes the 35 additional MR pairs across three clinical domains added in the expanded analysis. Together with the 26 MS/AD pairs in Table 12, these constitute the 61-pair catalog used in §4.1.

The pattern is consistent across all three new domains: transport pairs produce  $Q < 1$  with  $I^2 = 0\%$ , non-transport pairs with adequate power produce  $Q > 7$ , and all nine misclassifications (marked †) are false negatives with power  $< 0.38$ . The cardiometabolic domain provides particularly clean results because ancestry-stratified GWAS for lipids, blood pressure, and metabolic traits have large sample sizes, producing well-separated effect estimates.

## G bnlearn curvature benchmark details

The four benchmark DAGs used in the curvature validation are standard causal graphs from the bnlearn repository (Scutari, 2010): Asia (Lauritzen and Spiegelhalter, 1988), Sachs (Sachs et al., 2005), Insurance (Binder et al., 1997), and Alarm (Beinlich et al., 1989). For each graph, we generate linear SEM data ( $X_j = \sum_{i \in \text{pa}(j)} \beta_{ij} X_i + \epsilon_j$ ,  $\beta \sim \text{Uniform}(0.3, 0.8)$  with random sign,  $\epsilon \sim \mathcal{N}(0, 1)$ ,  $n = 1,000$ ) and learn a DAG via a PC-style algorithm using partial correlation CI tests with Fisher  $z$ -transform, conditioning sets up to size 2, and v-structure orientation. Edge labels are assigned by skeleton match: a learned edge is TP if it exists (in either direction) in the true DAG, FP otherwise. AUROC is computed via the Mann–Whitney  $U$  statistic across 50 replicates per (graph,  $\alpha$ ) combination.

The key finding is that the TP/FP discrimination direction depends on graph topology. On random DAGs (main text), true edges tend to connect lower-degree nodes, so Forman–Ricci ( $\kappa_F = 4 - d(u) - d(v)$ ) assigns higher curvature to TP edges (AUROC  $> 0.50$ ). On the benchmark DAGs, hub nodes (e.g., PKC and PKA in Sachs, with 5+ outgoing causal edges each) carry many true edges, so TP edges have *higher* endpoint degrees and Forman–Ricci assigns them *lower* curvature (AUROC  $< 0.50$ ). Betweenness centrality, which correlates with both degree and structural importance, remains the most robust discriminator across both random and structured topologies.

Table 21: MS MR pair sources. Each pair uses 3–4 strata from published multi-ancestry or multi-cohort genetic studies.

Pair	Stratum	$\hat{\beta}$	SE	Source
HLA-DRB1 → MS risk	EUR (IMSGC)	0.93	0.03	IMSGC 2019 <i>Science</i> , EUR meta ( $n > 47k$ cases)
	AFR	0.55	0.12	Isobe 2015 <i>Neurol. Genet.</i> ( $n \approx 1,500$ )
	EAS	0.40	0.15	Yoshimura 2012, Japanese MS HLA
	HISP	0.72	0.10	Isobe 2015, Hispanic MS HLA
Latitude/vitD gradient	EUR high lat	0.45	0.05	Simpson 2011 <i>Neurology</i> , high-latitude
	EUR low lat	0.20	0.06	Simpson 2011, Mediterranean
	EAS	0.10	0.08	Wallin 2019 <i>Lancet Neurol.</i>
	AFR	0.05	0.10	Wallin 2019
Sex ratio (F:M)	EUR recent	1.10	0.04	Walton 2020 <i>Lancet Neurol.</i> (F:M $\approx$ 3:1)
	EUR historical	0.69	0.06	Walton 2020 (F:M $\approx$ 2:1)
	EAS (Japan)	0.85	0.08	Osoegawa 2009, Japanese MS
	AFR	1.20	0.10	Langer-Gould 2013, Black MS
EBV → MS (causal)	EUR (Bjornevik)	1.50	0.15	Bjornevik 2022 <i>Science</i> , US military
	EUR (UKB MR)	1.35	0.20	Harroud 2024, MR EBV→MS
	Multi-ethnic	1.60	0.25	Ascherio 2012, seroconversion meta
Smoking → MS (MR)	EUR (IMSGC)	0.15	0.04	Hedstrom 2016 + Harroud 2021
	EUR (UKB)	0.12	0.05	UKB instruments + IMSGC outcome
	EUR (Scand.)	0.18	0.06	Hedstrom 2013, Scandinavian
Vitamin D → MS (MR)	EUR (IMSGC)	-0.25	0.06	Mokry 2015 <i>PLoS Med</i>
	EUR (repl.)	-0.20	0.08	Rhead 2016 <i>Neurology</i>
	EUR (Scand.)	-0.30	0.10	Gianfrancesco 2017
BMI → MS (MR)	EUR (IMSGC)	0.18	0.05	Mokry 2016 <i>PLoS Med</i>
	EUR (UKB)	0.15	0.06	Harroud 2020, BMI MR
	EUR (Scand.)	0.20	0.07	Gianfrancesco 2017

Table 22: AD MR pair sources. Non-transport pairs use multi-ancestry strata; transport pairs use multi-cohort EUR replication. Approximate estimates (marked †) are calibrated from related studies. APOE4 ORs predate the Belloy meta-analysis (Belloy et al., 2023) ( $n = 68,756$ ), which reports updated ancestry-stratified dosage ORs; tau spreading coefficients predate the Smith et al. longitudinal analysis (Smith et al., 2020).

Pair	Stratum	$\hat{\beta}$	SE	Source
APOE4 → AD risk	EUR	1.241	0.028	Belloy 2023 <i>JAMA Neurol</i> , $\varepsilon_3/\varepsilon_4$ ( $n = 34,021$ )
	AFR	0.779	0.069	Belloy 2023, Black ( $n = 7,145$ )
	EAS	1.513	0.066	Belloy 2023, East Asian ( $n = 21,852$ )
	HISP	0.642	0.071	Belloy 2023, Hispanic ( $n = 5,738$ )
Sex × tau spread	Entorhinal, F	0.42	0.06	Buckley 2019 <i>Ann Neurol</i> , HABS/ADNI women
	Entorhinal, M	0.18	0.07	Buckley 2019, men
	Temporal, F	0.38	0.05	Buckley 2019, women temporal
	Temporal, M	0.15	0.06	Buckley 2019, men temporal
BMI → AD (MR)	EUR (midlife)	-0.18	0.05	Nordestgaard 2017, EUR midlife BMI MR
	EUR (late-life)	0.05	0.06	Nordestgaard 2017, late-life (null)
	EAS†	0.12	0.08	Cross-ancestry, EAS direction reversed
	AFR†	-0.08	0.10	Reitz 2020, AFR BMI-AD MR
CRP → AD (MR)	EUR (UKB)	0.010	0.030	Rasmussen 2019 <i>Brain</i> , MR using UKB
	EUR (IGAP)	-0.020	0.035	Zheng 2020 <i>BMC Med</i> , MR using IGAP
	EUR (Bellenguez)†	0.015	0.028	CRP MR using Bellenguez 2022 outcome
TREM2 R47H → AD	EUR (meta)	1.131	0.12	Guerreiro 2013 + Jonsson 2013, EUR meta
	EUR (repl.)	0.956	0.18	Sims 2017 <i>Nat Genet</i> , EUR replication
	EUR (UKB)†	1.030	0.20	UK Biobank AD-by-proxy GWAS
Education → AD (MR)	EUR (Kunkle)	-0.095	0.040	Larsson 2017 <i>Ann Neurol</i> , education MR
	EUR (Bellenguez)†	-0.110	0.035	Updated MR with Bellenguez 2022
	EUR (UKB proxy)†	-0.080	0.045	UKB AD-by-proxy outcome
Alcohol → AD (MR)	EUR (UKB)	0.020	0.040	Larsson 2017 <i>Ann Neurol</i> , alcohol MR
	EUR (IGAP)†	-0.010	0.050	Alcohol MR using IGAP outcome
	EUR (Bellenguez)†	0.030	0.035	Alcohol MR using Bellenguez 2022
Ancestry × threshold	EUR (florbetapir)	0.89	0.03	Landau 2013, ADNI florbetapir sensitivity
	EUR (florbetaben)	0.82	0.04	Bullich 2017, florbetaben EUR

Table 23: Expanded MR pair catalog: 35 pairs across cardiometabolic disease, cancer, and psychiatric disorders. All pairs use 3–4 ancestry or cohort strata from published multi-ancestry GWAS or MR studies.

	<b>Domain</b>	<b>Pair</b>	<b>Expected</b>	<b><i>k</i></b>	<b>Correct</b>
Cardiometabolic		LDL → CAD	transport	4	Yes
		SBP → stroke	transport	4	Yes
		HDL → CAD	transport	4	Yes
		Triglycerides → CAD	transport	3	Yes
		Lp(a) → CAD	transport	3	Yes
		Smoking → CAD	transport	3	Yes
		T2D → CAD	transport	4	Yes
		SBP → CAD	transport	4	Yes
		Fasting glucose → T2D	transport	3	Yes
		LDL → stroke	transport	3	Yes
		IL-6R → CAD	transport	3	Yes
		BMI → asthma	transport	4	Yes
		BMI → T2D	non-tr	4	Yes
		Alcohol → CAD	non-tr	4	Yes
		BMI → CAD	non-tr	4	Yes
		WHR → T2D	non-tr	4	Yes
		Urate → gout	non-tr	4	Yes
		Alcohol → cirrhosis	non-tr	3	Yes
Cancer		Smoking → lung cancer	transport	4	Yes
		Alcohol → breast cancer	transport	3	Yes
		BMI → colorectal cancer	transport	4	Yes
		Height → prostate cancer	transport	3	Yes
		Smoking → bladder cancer	transport	3	Yes
		BMI → endometrial cancer	transport	3	Yes
		Smoking → COPD	transport	3	Yes
		BMI → kidney cancer	transport	3	Yes
		Telomere length → cancer	transport	3	Yes
		BMI → breast cancer	non-tr	4	Yes
		Insulin res. → pancreatic <sup>†</sup>	non-tr	4	No
Psych.		Cannabis → schizophrenia	transport	3	Yes
		Smoking → schizophrenia	transport	3	Yes
		CRP → depression	transport	3	Yes
		BMI → bipolar	transport	3	Yes
		BMI → depression <sup>†</sup>	non-tr	4	No
		Education → schiz. <sup>†</sup>	non-tr	4	No