1 **Revealing a coherent cell state landscape across single cell datasets with**

2 **CONCORD**

3 Qin Zhu[1]*, Zuzhi Jiang[2], Matt Thomson[3], Zev Gartner[1,4,5]*

4

5 [1] University of California San Francisco, Department of Pharmaceutical Chemistry, San

6 Francisco, CA 94158, USA

7 [2] Tetrad Graduate Program, University of California San Francisco, San Francisco, CA

8 94158, USA.

9 [3] Division of Biology and Biological Engineering, California Institute of Technology,

10 Pasadena, CA 94158, USA

11 [4] Chan Zuckerberg Biohub, University of California San Francisco, San Francisco, CA

12 94158, USA

13 [5] Center for Cellular Construction, University of California San Francisco, CA 94158,

14 USA

15

16 *Correspondence to Q. Zhu at qin.zhu@ucsf.edu and Z. J. Gartner at

17 zev.gartner@ucsf.edu

18

19 **Author information**

20 Qin Zhu, qin.zhu@ucsf.edu

21 Zuzhi Jiang, zuzhi.Jiang@ucsf.edu

22 Matt Thomson, mthomson@caltech.edu

23 Zev J. Gartner, zev.gartner@ucsf.edu

24

25

## Abstract

Resolving the intricate structure of the cellular state landscape from single-cell RNA sequencing (scRNAseq) experiments remains an outstanding challenge, compounded by technical noise and systematic discrepancies—often referred to as batch effects—across experimental systems and replicate. To address this, we introduce CONCORD (COntrastive learNing for Cross-dOmain Reconciliation and Discovery), a self-supervised contrastive learning framework designed for robust dimensionality reduction and data integration in single-cell analysis. The core innovation of CONCORD lies in its probabilistic, dataset- and neighborhood-aware sampling strategy, which enhances contrastive learning by simultaneously improving the resolution of cell states and mitigating batch artifacts. Operated in a fully unsupervised manner, CONCORD generates denoised cell encodings that faithfully preserve key biological structures, from fine-grained distinctions among closely related cell states to large-scale topological organizations. The resulting high-resolution cell atlas seamlessly integrates data across experimental batches, technologies, and species. Additionally, CONCORD's latent space capture biologically meaningful gene programs, enabling the exploration of regulatory mechanisms underlying cell state transitions and subpopulation heterogeneity. We demonstrate the utility of CONCORD on a range of topological structures and biological contexts, underscoring its potential to extract meaningful insights from both existing and future single-cell datasets.

## Introduction

Cells express thousands of genes to perform specialized functions and maintain homeostasis. These gene expression patterns are governed by gene regulatory networks and cell–cell interactions, confining cell distributions to a low-dimensional state space (or 'state landscape') within the high-dimensional gene expression space[1,2]. Waddington's landscape is often invoked as a conceptual model of the state landscape, illustrating how cells traverse branching developmental trajectories toward distinct fates[3]. Modern single-cell technologies such as scRNA-seq capture snapshots of these processes, revealing diverse topological structures - including clusters, bifurcations, convergences and loops—indicating more nuanced structures than that depicted in Waddington's original diagram[4]. Accurately capturing these structures, which are typically unknown a priori, is therefore essential for elucidating developmental processes, homeostatic regulation, and disease progression.

Dimensionality reduction, or representation learning in machine learning, is a common approach for uncovering such structures. By projecting high-dimensional data into a lower-dimensional space, key structural patterns become more tractable to visualize and analyze. However, standard techniques—such as principal component analysis, non-negative matrix factorization[5], and factor analysis[6]—tend to overemphasize broad cell type separation, overlook subtle states, and can confound processes like differentiation with cell cycle progression. Compounding these challenges is the continuous nature of cell state changes, yet scRNA-seq profiles typically sample only discrete snapshots. Researchers often integrate multiple datasets capturing distinct developmental stages, biological replicates, perturbations, or other conditions to assemble a more comprehensive view of the underlying state landscape. Yet, such integration efforts are often impeded by technical and biological batch effects, which can obscure genuine biological signals. Although an array of batch-correction tools—such as Harmony[7], Scanorama[8], Seurat[9], scVI[10], LIGER[11] and

3

78  MNN[12] —have been developed, they can distort the biological structure by over-

79  or under-correcting batch effects[13], and many face scalability issues when

80  applied to massive atlas-level datasets.

81

82  Recently, contrastive learning—a technique proven effective in image and natural

83  language processing[14-16]—has gained attention in single-cell genomics[17-23]. This

84  approach learns meaningful representations by contrasting similar (positive) cells

85  against dissimilar (negative) ones within mini-batches—small, randomly sampled

86  subsets of the dataset that are iteratively contrasted during training. By building

87  an understanding of what makes cells similar and different, contrastive learning

88  extracts features underly distinct cellular states, thereby improving clustering and

89  annotation[17-19,21-23]. Some methods adopt a supervised strategy, using known cell

90  type labels to define positive and negative pairs[22,23]. However, this approach

91  requires substantial manual annotation effort and limits generalization to novel

92  cell states or continuous trajectories. Others employ self-supervised strategies,

93  treating each cell in a mini-batch as a negative for all others[17-20]. Yet, with

94  negatives uniformly drawn from the global distribution, these methods often

95  emphasize global differences - such as major cell type clusters - while

96  underutilizing the fine resolution of single-cell data to detect subtle cellular states.

97  Although recent work suggests that mini-batch composition heavily influences

98  contrastive outcomes[24], existing single-cell approaches have yet to fully explore

99  the implications and generalizability of this insight.

100

101  Applying contrastive learning to multiple datasets presents additional challenges:

102  naively contrasting cells across different batches can amplify batch-specific

103  artifacts rather than isolating biologically relevant variation. Various techniques -

104  such as generative adversarial networks (GANs)[20,25,26], unsupervised domain

105  adaptation via backpropagation[27], and conditional variational autoencoders

106  (CVAEs)[28] - have been employed to mitigate batch effects. However, in a

107  contrastive learning setting, maximizing differences between dissimilar cells

108    conflicts with the goal of minimizing dataset-specific differences, frequently

109    resulting in incomplete removal of batch effects. This dilemma raises a key

110    question: can contrastive learning fully capture cellular diversity while minimizing

111    batch artifacts?

112

113    In this work, we introduce CONCORD (COntrastive learNing for Cross-dOmain

114    Reconciliation and Discovery), a novel contrastive learning framework designed

115    to denoise and reveal the detailed structure of the cellular state landscape across

116    single or multiple datasets. Central to CONCORD is a probabilistic,

117    neighborhood-aware and dataset-aware sampling strategy that harnesses mini-

118    batch composition to achieve two critical goals: enhancing the resolution of cell

119    states in the latent space, and removing batch-associated biases and noise.

120    Neighborhood-aware sampling prompts the model to contrast nearby cells,

121    thereby capturing subtle distinctions among closely related states. Meanwhile,

122    dataset-aware sampling ensures comparisons occur within the same dataset,

123    prioritizing biological variation while minimizing dataset-specific artifacts.

124    Crucially, these strategies are unified within a single probabilistic sampler,

125    enabling contrastive learning to simultaneously improve resolution and mitigate

126    batch effects. Using simulated and real datasets, we demonstrate that

127    CONCORD produces high-resolution, denoised encodings that robustly capture

128    diverse structures—including loops, trajectories, trees, and specialized cell

129    states—reflecting bona fide biological processes even when the data originate

130    from multiple technologies, time points, or species.

131

132    **Results**

133    **The CONCORD framework**

134    Analysis of single-cell sequencing data suggest that gene expression is not

135    randomly sampled; rather, the mechanism of gene regulation imposes strong

136    constraints, producing dynamically changing gene co-expression patterns

5

137   reflected as intricate structures in the low dimensional embedding of cells[1,2,4]. For

138   example, adult cells at homeostasis often form discrete clusters corresponding to

139   stable cell types or states, with adjacent clusters representing closely related cell

140   states (Figure 1A, left). In contrast, developmental or pathological contexts—

141   such as early embryogenesis, tissue repair, or tumorigenesis—tend to exhibit

142   trajectories branching from progenitor cells to multiple terminal fates, with semi-

143   stable transitional states forming denser clusters (Figure 1A, middle). Loop

144   structures are also frequently observed[4,29], typically associated with cyclic gene

145   expression programs, such as those governing cell cycle progression[26] (Figure

146   1A, right). Despite these rich patterns, standard dimensionality reduction

147   methods like principal component analysis (PCA) or nonnegative matrix

148   factorization (NMF) capture only partial representations of the cell state

149   landscape, either oversimplifying complex structures or disproportionately

150   emphasizing certain features while obscuring others.

151

152   We hypothesized that a latent representation method capable of intrinsically

153   modeling relationships among cells at the level of gene co-expression programs

154   could yield a more comprehensive view of the state landscape. Recent evidence

155   suggests that self-supervised contrastive learning, particularly SimCLR-type

156   approaches[14], significantly improves clustering and cell-type classification

157   performance[18-20]. In this framework, positive pairs are typically generated by

158   perturbing a given cell—such as by randomly masking gene expression values—

159   to create two augmented versions. Contrastive learning then operates through

160   two complementary forces: differentiation and alignment. By contrasting each cell

161   against all others in the mini-batch, the model identifies features that differentiate

162   between distinct cell states. Simultaneously, aligning the two masked versions of

163   the same cell encourages the model to capture gene co-expression patterns,

164   rather than relying on the expression of individual genes[30]. Theoretical results

165   demonstrate that contrastive learning with random masking and ReLU networks

166   effectively captures sparse signals while suppressing noise[30] (see Methods). In

167   the context of gene expression, these sparse signals correspond to gene co-

168  expression programs, where sets of co-expressed genes define functional cell

169  states[6,31,32]. We incorporated similar conditions—LeakyReLU activation and

170  random masking—into our model architecture (see Methods, Supplemental

171  Figure 1) and observed an encoding that effectively captures gene co-expression

172  programs with suppressed universal features and noise (Figure 1B). By

173  prioritizing sets of co-expressed genes over individual gene expression

174  fluctuations, this representation is also inherently more robust to dropouts – a

175  common artifact in single-cell RNA-seq[33].

176

177  In standard contrastive learning, training proceeds iteratively: in each step, a

178  mini-batch is drawn from the dataset, cells are augmented, encodings are

179  computed, and a contrastive loss is applied based on similarities between

180  augmentations and differences among samples. This cycle repeats over multiple

181  iterations, gradually refining the latent encodings. Since contrastive objectives

182  are computed within each mini-batch[24], the sampling strategy – which dictates

183  mini-batch composition - plays a pivotal role in shaping the learned

184  representation. Most methods rely on uniform sampling across the entire dataset,

185  leading to two key limitations. First, uniform sampling emphasizes broad

186  differences, such as major cell types, while underrepresenting rare

187  subpopulations, leading to poor resolution of fine-scale cellular states (Figure

188  1B). Second, mixing cells from different datasets within the same mini-batch can

189  amplify dataset-specific differences, inadvertently encoding batch effects rather

190  than isolating biologically meaningful variation.

191

192  To address the first issue, we developed a neighborhood-aware sampler, inspired

193  by k-nearest-neighbor (kNN) based sampling[24]. In this approach, cells are

194  sampled probabilistically from both the global distribution and local

195  neighborhoods (Figure 1B). Local sampling, guided by a coarse graph

196  approximation of the cellular state landscape, compels the model to contrast cells

197  against their neighbors, allowing it to capture subtle differences between closely

198 related states. Meanwhile, global sampling preserves a broad perspective of

199 major cell types, ensuring the model robustly encodes large-scale distinctions. By

200 iteratively presenting the model with local neighborhoods (e.g., T cells in one

201 mini-batch, epithelial cells in another) alongside the global distribution, the model

202 allocates capacity to represent both large-scale distinctions and nuanced local

203 details, leading to improved resolution of biologically structures in the learned

204 latent space (Figure 1B).

205

206 When applied to a single dataset, contrastive learning effectively distills biological

207 variation into a meaningful latent representation (Figure 1C). However, with

208 uniform sampling across multiple datasets, both biological and dataset-specific

209 variations are encoded, leading to latent spaces that separate by dataset as well

210 as cell type (Figure 1D). To mitigate these batch effects, we additionally

211 developed a dataset-aware sampling approach. This method leverages the fact

212 that contrastive learning primarily captures biologically variable signals within

213 each mini-batch, while suppressing global, non-informative features and noises[30]

214 (Figure 1B). By restricting mini-batches to a single dataset, the model learns only

215 biological variations at each training step, mirroring the effect of applying

216 contrastive learning within a single dataset. Meanwhile, batch-associated biases

217 and noise are progressively reduced through mini-batch shuffling, which

218 randomly interleaves mini-batches from different datasets. If the contrastive loss

219 induces dataset-specific encodings in one mini-batch, subsequent mini-batches

220 from other datasets disrupt and overwrite these encodings due to their

221 inconsistency across datasets. As a result, only biologically meaningful signals,

222 such as gene co-expression patterns, persist across training iterations, leading to

223 a latent space that accurately reflects cell type and state distributions with

224 significantly reduced batch effects (Figure 1E).  In cases where datasets have

225 minimal or no overlap, a "leaky" dataset-aware sampler enables soft alignment

226 without imposing artificial harmonization, allowing for a flexible integration

227 scheme that generalizes to both fully and partially overlapping datasets (Figure

228 1F). Unlike traditional batch-correction approaches used in machine learning

229  models—such as conditional variational autoencoders (CVAEs)[28], unsupervised

230  domain adaptation by backpropagation[27], and generative adversarial networks

231  (GANs)[20,25,26] —which often struggle in contrastive learning due to competing

232  objectives, CONCORD directly leverages the contrastive objective and training

233  process to learn a batch-effect-mitigated latent representation.

234

235  Both the neighborhood-aware and dataset-aware samplers follow a unified

236  design principle: probabilistically structuring mini-batches to balance global

237  biological variations with local and dataset-specific differences. In the

238  neighborhood-aware sampler, most cells are selected to maintain the global

239  distribution of cell types, ensuring stable representation of large-scale biological

240  variation, while a smaller proportion is drawn from local neighborhoods to

241  enhance fine-grained resolution. Similarly, in the dataset-aware sampler, mini-

242  batches primarily consist of cells from a single dataset to minimize dataset-

243  specific biases, with a smaller subset incorporating cross-dataset samples to

244  retain biologically meaningful variation. By integrating both samplers, we

245  establish a joint probabilistic sampling framework, where the likelihood of

246  selecting a given cell reflects the combined probabilities of dataset-aware and

247  neighborhood-aware sampling (Figure 1F). This generalized sampling strategy

248  enables both robust dataset integration and improved resolution within a self-

249  supervised contrastive learning framework, forming the foundation of the

250  CONCORD model (Supplemental Figure 1). Notably, we found that a simple two-

251  layer encoder with Leaky ReLU activation effectively captures rich biological

252  structures without requiring additional layers or a decoder, significantly reducing

253  training data requirements and improving robustness. This minimalistic

254  architecture also enhances interpretability, as a Leaky ReLU activation after the

255  first linear layer enables the second layer to encode context-dependent gene

256  expression programs, effectively compressing complex cell states into a

257  moderate number of latent neurons.

258

**CONCORD learns denoised latent representations that preserve underlying structures**

Recovering biologically meaningful insights from single-cell data relies on preserving key features of the gene expression state space, including both its geometric organization and topological structure. To assess whether CONCORD meets this requirement, we evaluated its performance on both simulated and real-world single-cell datasets. While existing simulation packages (e.g., splatter[34]) generate discrete cell-type clusters, they fail to recapitulate the diverse structures observed in real biological data, including branching trajectories, loops, and multi-scale hierarchies. To address this limitation, we developed a custom simulation workflow capable of producing a wide range of structures that closely resemble real biological data, while allowing flexible control over noise distributions and batch effects (Figure 2A).

In parallel, we designed an evaluation pipeline that integrates both geometric and topological metrics to quantitatively assess the quality of latent representations (Figure 2B). Traditional benchmarking frameworks, such as scIB[35], primarily focus on biological label preservation and batch correction, but provide less insight into whether a representation retains the intrinsic geometrical relationships and topological structures of the original data. To fill this gap, we incorporated geometric metrics, including trustworthiness[36] and distance correlation, alongside topological data analysis (TDA) methods based on persistent homology and Betti numbers (Figure 2B). Trustworthiness quantifies how well local neighborhood structures in high-dimensional space are preserved in the latent representation, with high scores indicating that nearby cells in the original data remain neighbors in the learned embedding. Persistent homology, in contrast, provides a more global perspective by tracking the emergence and persistence of topological features—such as connected components (i.e. clusters; Betti-0), loops (Betti-1), and enclosed voids (Betti-2)—across varying distance scales. These features are visualized in persistence diagrams and Betti

10

289     curves, where stable structures appear as long-lived features in the persistence

290     diagram and extended plateaus in the Betti curve, whereas transient, noise-

291     induced features vanish quickly.

292

293     We first evaluated CONCORD on a simple simulated dataset consisting of three

294     well-separated clusters corrupted by cluster-specific Gaussian noise (Figure 2C,

295     Supplemental Figure 2A). We compared CONCORD against a diverse set of

296     dimensionality reduction and embedding methods, including diffusion map, NMF,

297     Factor Analysis (FA), FastICA, Latent Dirichlet Allocation (LDA), ZIFA, scVI, and

298     PHATE. While visualizations with the UMAP algorithm[37] provide a qualitative

299     comparison (Figure 2C), all benchmarking statistics were computed directly on

300     the latent space, ensuring that evaluations reflect the true structure learned by

301     each model rather than artifacts introduced by UMAP. Many methods, including

302     NMF, FA, FastICA, ZIFA, and scVI, struggled to fully separate the clusters, while

303     others, such as PHATE, introduced spurious trajectory-like structures (Figure

304     2C). By contrast, CONCORD cleanly distinguished the three clusters and

305     produced a denoised distance matrix that closely matched the ground truth

306     (Figure 2C). Persistent homology further confirmed these results: CONCORD's

307     persistence diagram correctly identified the expected three-cluster topology

308     (Betti-0 = 3), showing a stable plateau in the Betti curve that closely matched the

309     noise-free dataset.

310

311     For more complex structures, such as a self-connecting trajectory with three

312     loops and multiple branching points, CONCORD was the only method that

313     accurately recovered the full underlying structure (Figure 2D, Supplemental Fig.

314     2B). Other methods either collapsed trajectories into clusters (e.g., DiffusionMap,

315     NMF, FactorAnalysis, FastICA, LDA) or failed to generate a stable persistence

316     diagram, likely due to excessive noise retention (e.g., ZIFA, scVI). While PHATE

317     produced a relatively smooth latent space (Figure 2D), its Betti curve analysis

318    revealed only a single persistent loop rather than the expected three, indicating

319    an incomplete preservation of topological structure.

320

321    We summarized these findings using geometric and topological metrics and

322    found that CONCORD consistently outperformed alternative methods across a

323    broad range of simulated structures (Figure 2E, 2F). Importantly, despite its

324    denoising capability, CONCORD preserves relative noise levels in the latent

325    space, as demonstrated by the strong correlation between latent variance and

326    input variance in the cluster simulation (Figure 2E). This is particularly important

327    because biological noise—such as transcriptional variability—often carries

328    biologically relevant information, distinguishing dynamic cellular states from

329    stable ones. Furthermore, CONCORD consistently maintains high

330    trustworthiness scores across neighborhood sizes from 10 to 100, underscoring

331    its ability to preserve local structure over a significant radius (Supplemental

332    Figure 2C, 2D). In contrast, many alternative methods exhibit sharp declines in

333    trustworthiness, indicating a loss of fine-scale geometric relationships in the

334    latent space.

335

336    To investigate how neighborhood-aware sampling influences performance, we

337    simulated a three-level branching tree that mimics hierarchical tissue-cell-type-

338    cell-state relationships (Figure 2G, Supplemental Figure 2E). Without the

339    neighborhood-aware sampler, the model failed to resolve sub-branches; by

340    contrast, higher local sampling uncovered more subtle differences (Supplemental

341    Figure 2E, 2F), leading to substantially improved resolution in the latent

342    representation (Figure 2G). However, excessive local enrichment had

343    drawbacks. When intra-kNN sampling probability exceeded 0.6, global

344    distinctions at tissue and cell-type levels were suppressed (Figure 2G,

345    Supplemental Figure 2F), consistent with our earlier findings that contrastive

346    learning harmonizes inconsistent global signals. Based on these results, we

347 recommend an intra-kNN sampling probability below 0.5 to maintain a balance

348 between fine-grained resolution and global structure preservation.

349

350 **CONCORD learns a coherent, batch-effect-mitigated latent representation**

351 Batch effects often manifest as global signals that impact all cells within a dataset

352 but differ across datasets. When mini-batches are enriched for cells from a single

353 dataset, these dataset-specific signals rapidly diminish during training due to their

354 inconsistent nature (Figure 1E, Figure 3A). This batch-correction mechanism

355 fundamentally differs from conventional methods, which typically assume a

356 specific model of batch distortion and force the alignment between datasets. By

357 contrast, CONCORD prioritizes the learning of biologically informative and

358 coherent variation over inconsistent batch signals, making minimal assumptions

359 about the nature of batch effects. As a result, it preserves biological structures

360 more effectively while mitigating technical artifacts.

361

362 To evaluate this, we tested CONCORD on a simulated dataset containing five

363 clusters corrupted by both batch effects and noise. Among all methods tested,

364 CONCORD was the only approach that accurately identified all clusters;

365 alternative methods struggled to separate closely related cell types or introduced

366 alignment artifacts, such as the "ring" structure produced by Scanorama[8].

367 Notably, applying a standard contrastive learning sampler without dataset-aware

368 sampling resulted in pronounced batch effects (Figure 3A). While significantly

369 reducing noise, CONCORD preserved cluster-specific variance differences rather

370 than over-smoothing them—a critical aspect often overlooked by batch-correction

371 algorithms (Supplemental Figure 3A). Moreover, CONCORD demonstrated

372 consistent denoising and batch-effect correction performance across different

373 noise distributions (e.g. Gaussian, negative-binomial, Poisson) and batch-effect

374 models (e.g. batch-specific gene expression differences, sequencing depth

375 variability, and dropout rate differences), underscoring its independence from

376 specific noise and batch effect assumptions (see Methods; data not shown).

377

Beyond aligning discrete clusters, a major challenge in single-cell data integration is constructing a unified representation from datasets capturing related but distinct conditions—such as different developmental stages, tissues, species, or perturbations. Many integration techniques struggle in cases where datasets exhibit limited overlap in cell states, particularly those relying on matched clusters (e.g., Harmony[7]) or mutual nearest neighbors (e.g., Seurat[9], Scanorama[8], MNN[12]). Aligning datasets solely based on overlapping states can also distort the structure when these regions are sparse or imbalanced. As the number of datasets increases, these limitations become more pronounced, often leading to fragmented or warped embeddings.

388

To systematically assess these challenges, we simulated batch effects across a range of structures: clusters (Supplemental Figure 3B), trajectories (Figure 3B, Supplemental Figure 3C), loops (Figure 3C, Supplemental Figure 3D), and trees (Figure 3D, Supplemental Figure 3E). On a trajectory simulation where batches were fully overlapping, most methods achieved some degree of alignment, though several exhibited partial or suboptimal correction (e.g., scVI, LIGER). However, as overlapping regions were progressively reduced, most methods deteriorated. Some failed to align batches entirely (e.g., LIGER, Harmony), while others introduced spurious structures, such as Scanorama producing an artificial loop instead of a linear trajectory, or scVI collapsing trajectories with partial overlap into a single cluster. Across all conditions, CONCORD consistently recovered the underlying structure with significantly reduced noise, even when shared cell states were minimal or absent (Figure 3A–D, Supplemental Figure 3B–E, Supplemental Table 1).

403

CONCORD's robustness likely arises from its focus on learning biologically meaningful gene co-expression programs (see Methods, Figure 1B, Supplemental Figure 2A, 2B, 2E) rather than forcing direct cell-state alignment

14

407 across batches—a fundamental distinction from many existing batch-correction

408 approaches. By encoding gene co-expression programs, CONCORD positions

409 cells with similar transcriptomic states together in the latent space, eliminating

410 the need for explicit reference points. This property allows CONCORD to better

411 preserve biological signals, as reflected in its high biological label conservation

412 scores in the scIB benchmarking metrics (Figure 3E, 3F). However, because

413 CONCORD does not explicitly merge batches, it achieves high but not the

414 highest batch-correction scores. Interestingly, CONCORD more effectively

415 preserves local distance relationships than global ones, leading to high

416 trustworthiness scores but lower global distance correlations, which affects its

417 overall geometric score ranking (Figure 3F, Supplemental Table 1). This outcome

418 reflects an intrinsic trade-off in manifold learning, where maintaining local

419 neighborhood structure is often favored over enforcing global distance

420 relationships[38,39]. Nonetheless, CONCORD consistently ranks among the top-

421 performing methods for topological structure preservation and biological label

422 conservation, as well as in overall rankings (Figure 3F). These results indicate

423 that CONCORD effectively mitigates batch effects while preserving the intrinsic

424 structure of the data. Given the complexity of real single-cell datasets, these

425 findings suggest that CONCORD provides a reliable, generalizable solution for

426 both dimensionality reduction and batch-effect mitigation, even in cases with

427 uncertain underlying structures or limited batch overlap.

428

429 **CONCORD aligns whole-organism developmental atlases and resolves**

430 **high-resolution lineage trajectories**

431 To evaluate CONCORD's ability to capture biologically meaningful structures

432 across datasets generated by different technologies, we benchmarked it against

433 popular batch correction methods using lung and pancreas single-cell atlases[35]

434 (Supplemental Figure 4A, B). While CONCORD excelled at identifying discrete

435 cell-type clusters, these datasets and their annotations do not reflect the

436 hierarchical organization of biological systems and lack the continuous or

branching trajectories characteristic of developmental and pathological processes. A rigorous assessment of integration methods requires datasets with diverse structures and ground-truth lineage annotations, enabling a fine-scale evaluation of both dataset integration and biological structure conservation.

Caenorhabditis elegans (*C. elegans*) embryogenesis provides an ideal benchmark, as it follows a well-characterized, nearly invariant lineage tree from the fertilized egg to the 558 cells present at hatching[40]. This lineage tree is highly conserved in closely related species, such as *C. briggsae*[41]. Packer et al. previously generated a lineage-resolved single-cell atlas of *C. elegans* embryogenesis[42], which was recently expanded by Large et al. to include over 200,000 *C. elegans* cells and 190,000 *C. briggsae* cells[41]. The authors annotated progenitor and terminal cell fates through extensive subset-specific projections, clustering and fluorescent reporter imaging. These datasets serve as ideal resources for evaluating whether dimensionality reduction and integration algorithms can faithfully reconstruct and align developmental trajectories across species.

We first tested CONCORD on the *C. elegans* dataset[42] and observed that its UMAP embedding recapitulated known developmental trajectories (Supplemental Figure 5A). The effect of neighborhood-aware sampling mirrored trends observed in our simulations: moderate local enrichment improved resolution of subtle differences among neurons, while excessive local sampling disrupted the global structure (Supplemental Fig. 5A, B). This is because local sampling enhances mini-batches with underrepresented subpopulations, while adequate global sampling is necessary to capture broader variation (Supplemental Figure 5C).

Running CONCORD on the larger, cross-species dataset yielded a unified developmental atlas that closely matches original cell-type and lineage

466   annotations (Figure 4A). The embedding effectively separated broad cell classes

467   (Supplemental Figure 6A) and arranged cells along continuous trajectories from

468   early progenitors to terminal states in alignment with estimated embryo time

469   (Figure 4B). While UMAP is well-suited for visualizing complex latent

470   representations, the resulting latent representation proved too complex for two-

471   dimensional UMAP to disentangle fully— likely due to extensive lineage

472   bifurcations and convergences—necessitating a three-dimensional UMAP to

473   mitigate trajectory cross-overs (Figure 4A). Running CONCORD with or without a

474   decoder yielded similar UMAP embeddings (Figure 4B), and we present the with-

475   decoder version due to slightly less entangled trajectories. However, we strongly

476   encourage readers to explore the interactive 3D UMAP visualizations

477   (https://qinzhu.github.io/Concord_documentation/galleries/cbce_show/#__tabbed

478   _1_1), which offer a more detailed view of developmental trajectories than can be

479   captured in static 2D projections.

480

481   To systematically compare integration methods, we moved beyond UMAP-based

482   visual assessment and directly analyzed the 300-dimensional latent space

483   produced by each algorithm. Standard benchmarking tools, such as scIB[35], could

484   not scale to this dataset's size and complexity, requiring us to adapt previous

485   evaluation strategies. To assess species alignment, we computed species

486   composition within randomly sampled latent space neighborhoods, evaluating

487   whether each neighborhood's species fraction matched expectations, similar to

488   the kBET algorithm[43] (Figure 4C). Given the dataset's time-dependent sampling

489   biases, with *C. elegans* cells progressively increasing over time due to

490   experimental bias or true biological variation, we stratified the analysis by

491   developmental time bins. If species were well-aligned at a fine scale,

492   neighborhood species fractions should remain close to expected proportions with

493   minimal variation. Indeed, CONCORD maintained this pattern, whereas other

494   methods exhibited alignment failures: Scanorama produced neighborhoods

495   dominated by a single species (fractions of 0 or 1, also evident in the UMAP in

496   Figure 4B), scVI and LIGER showed excessive species fraction variability,

17

497  indicating incomplete species alignment, and Seurat and Harmony deviated

498  significantly from expected proportions, suggesting local density mismatches

499  between species.

500

501  Besides species alignment, another key benchmark is how well each method

502  preserves biological structure without sacrificing resolution. As seen in UMAP

503  visualizations, Harmony and scVI successfully aligned species but lost resolution,

504  obscuring subtypes and lineage trajectories (Figure 4B). Only Seurat and

505  CONCORD successfully recovered divergent terminal fates, and among them,

506  only CONCORD maintained continuous, fine-grained trajectories from

507  progenitors to terminal cell types (Figure 4B, D). Seurat's reciprocal PCA

508  approach captured later-stage variation but overshadowed subtle early lineage

509  differences, likely because later stages contain more cells and greater

510  transcriptional diversity. In contrast, CONCORD's neighborhood-aware sampling

511  preserved both early and late states, distributing representational capacity across

512  the entire developmental spectrum.

513

514  The rich lineage and cell-type annotations in this dataset enabled us to map the

515  *C. elegans* lineage tree onto UMAP embeddings, revealing how progenitor

516  lineages give rise to distinct terminal fates[40,44] (Supplemental 6B, C). Strikingly,

517  on many occasions, CONCORD's UMAP trajectories aligned with underlying

518  lineage subtrees. For example, ciliated amphid neurons ASE, ASJ, AUA—arising

519  from AB-derived progenitors—formed branching trajectories in CONCORD's

520  UMAP that closely mirrored the actual lineage sub-tree (Figure 4D). In contrast,

521  scVI and Seurat introduced large gaps between parent and daughter lineages or

522  obscured distinct terminal fates (Figure 4D). Beyond capturing lineage

523  relationships, CONCORD's latent activation patterns effectively distinguished

524  terminal cell states (Figure 4E). Hierarchical clustering of CONCORD's latent

525  space revealed two distinct ASE subpopulations, characterized by differential

526  expression of GCY receptors—gcy-5, gcy-19, gcy-22 in one group, and gcy-6,

18

527      gcy-7, gcy-14 in the other (Figure 4E). These receptors distinguish ASE-left

528      (ASEL) from ASE-right (ASER) neurons, which are morphologically symmetric

529      but functionally asymmetric in their salt-sensing responses[45,46].

530

531      To systematically assess the preservation of lineage structures in the latent

532      space, we repeated the lineage distance analysis from Packer et al.[42], correlating

533      latent space distances with lineage distances for AB lineages which give rise to

534      ~70% of terminal embryonic cells (Figure 4F). Consistent with the original study,

535      where transcriptome distances correlated with lineage distances, we found that

536      CONCORD's latent distances exhibited strong correlation with lineage distances,

537      even in early generations where transcriptomic differences were minimal. This

538      demonstrates CONCORD's ability to capture subtle, progressive molecular

539      changes underlying lineage bifurcations. Notably, CONCORD significantly

540      outperformed other integration methods in overall correlation (Figure 4G),

541      underscoring its potential for trajectory inference in developmental studies.

542

543      The improvement in resolving fine-scale structures became even more

544      pronounced when zooming into subsets of cells. In early embryonic cells,

545      Scanorama, Harmony, and scVI failed to fully align species or lost resolution,

546      whereas CONCORD revealed extensive lineage bifurcation patterns

547      (Supplemental Figure 6D). On muscle formation, CONCORD showed the MS, C,

548      and D lineages converge into sub-branches of body wall muscle, positioned from

549      the head (anterior) to the tail (posterior) in an orientation reflecting genuine

550      spatial gene expression gradients (Figure 4H, Supplemental Figure 6E).

551      CONCORD also captured rare lineage convergence events, such as

552      ABplp/ABprp-derived muscle integrating with MS-derived counterparts to form

553      intestinal muscles (mu_int) (Figure 4H). Pharyngeal development, involving

554      complex branching and convergence of AB- and MS-derived cells, was likewise

555      resolved by CONCORD (e.g., pm3–5 deriving from both AB and MS lineages,

556      and pm1–2, 6–8 specific to AB/MS lineage), whereas scVI and Seurat UMAPs

557 displayed fewer fine-grained details of such lineage convergence and divergence

558 (Figure 4I, Supplemental Figure 6F). Importantly, all of these zoom-in analyses

559 were performed directly on the global latent space learned by CONCORD,

560 without requiring subset-specific variable gene (VEG) selection and re-alignment

561 – steps that are often recommended for other methods.

562

563 Finally, CONCORD demonstrated superior efficiency and scalability on this

564 400,000-cell, 20-batch dataset (Figure 4J). It completed integration in ~30

565 minutes on an NVIDIA A100 GPU, outperforming LIGER, Seurat, and scVI, which

566 required several hours and significantly more memory. These results confirm that

567 CONCORD excels in scalability and accuracy, making it an ideal tool for

568 reconstructing complex developmental trajectories.

569

570 **CONCORD captures cell cycle and differentiation trajectories in mammalian**

571 **intestinal development**

572 Unlike *C. elegans*, where embryonic divisions are largely governed by maternally

573 supplied mRNAs[47], mammalian development involves producing many more

574 cells in each lineage, often coupling cell cycle activity with ongoing differentiation.

575 To evaluate whether CONCORD can recover these intertwined processes, we

576 applied it to a single-cell atlas of embryonic mouse intestinal development[48]. This

577 dataset encompasses multiple developmental stages, spatial segments, and

578 diverse cell types (Figure 5A). Moreover, it contains batches enriched for specific

579 populations—such as mesenchymal cells—providing a challenging scenario

580 where incomplete coverage across batches could hinder integration and

581 resolution of fine-grained cell states.

582

583 CONCORD not only integrated these datasets significantly faster than most other

584 methods (data not shown), but it also resolved finer substructures within each

585 cell type (Figure 5A, Supplemental Fig. 7A, Supplemental Figure 8). Batches

586    derived from different developmental stages and enrichment strategies merged

587    seamlessly into a cohesive cell atlas (Supplemental Figure 7A). Notably,

588    CONCORD revealed numerous loop-like patterns within each cell type,

589    corresponding to cell cycle progression (Figure 5A-D). Erythrocytes—known to

590    lack proliferative capacity—did not form such loops, matching established biology

591    (Supplemental Figure 8). As discussed earlier, because CONCORD can capture

592    complex structures, a standard 2D UMAP does not fully represents these loops

593    or breaks trajectories. We therefore focused our analysis on 3D UMAPs, (Figure

594    5A-D, Supplemental Figure 7B) and strongly recommend readers explore an

595    interactive visualization of CONCORD embedding

596    (https://qinzhu.github.io/Concord_documentation/galleries/huycke_show/).

597

598    In intestinal epithelial cells, CONCORD not only delineated detailed cell subtype

599    structures, including rare enteroendocrine cells (EECs), but also revealed two

600    parallel differentiation trajectories, each forming its own cell cycle loop (Figure

601    5B). Coloring the UMAP by zonation suggests that these trajectories correspond

602    to proximal and distal segment identities, reinforced by the expression of regional

603    markers such as *Bex1* and *Onecut2*[49](Figure 5B). While zonation in adult mouse

604    and human intestines has been documented[49], our data suggest that segment-

605    specific epithelial subsets arise as early as embryonic day 13.5 (E13.5), implying

606    that CONCORD can detect subtle positional signatures well before maturity.

607

608    In enteric nervous system (ENS) cells, CONCORD uncovered a clear cell cycle

609    loop alongside two bifurcating differentiation branches (Figure 5C). Morarach et

610    al.[50] previously characterized ENS progenitor states in mouse embryos, showing

611    progressive cell cycle phases and bifurcations marked by genes such as *Etv1*

612    and *Bnc2*. CONCORD successfully recapitulated these branching patterns,

613    appropriately anchoring them at a neuroblast population defined by *Cck*, which

614    arises from *Sox10*[+] progenitor cells. The convergence of the two branches

615    appears to be driven by a set of neuronal maturation genes that are broadly

616  expressed across both trajectories (Supplemental Figure 7C). Notably,

617  CONCORD was the only method that preserved both the cell cycle loop and the

618  bifurcation, whereas other methods introduced discontinuities or misplaced

619  bifurcation points (Supplemental Figure 8).

620

621  Mesenchymal cells comprise a major fraction of this dataset, yet prior studies

622  primarily focused on the *Pdgfra+* mesenchymal subset, which plays a key role in

623  villus formation[48]. In CONCORD's embedding, we not only recapitulated the

624  previously described *Pdgfra+* trajectory (Supplemental Figure 8), but also

625  uncovered highly heterogeneous substructures within the *Pdgfra−* mesenchymal

626  subset. This population, along with smooth muscle cells, formed four consecutive

627  cell cycle loops, each marked by distinct expression of *Ebf1*, *Slit2*, *Kit*, and *Acta2*,

628  with varying degrees of gradual transitions between loops (Figure 5D).

629  Interestingly, *Ebf1*, *Slit2* have been previously associated with mesenchymal cell

630  multipotency[51,52], while Kit marks interstitial cells of Cajal (ICC) and its

631  progenitors[53,54]. Unlike traditional approaches, where cell cycle effects often

632  obscure cell-state annotations, CONCORD preserves both cell cycle and

633  differentiation signals, allowing for the identification of these previously

634  unrecognized subpopulations.

635

636  The rich structures revealed by CONCORD is also evident in its latent space,

637  where each neuron is activated in one or more cell populations—ensuring no

638  latent capacity is wasted (Figure 5E). In contrast, the PCA results from Seurat's

639  integrated dataset compress variation into only a few principal components, while

640  scVI leaves nearly half of its latent dimensions minimally utilized. Furthermore,

641  CONCORD's latent space can be interrogated with gradient-based attribution

642  methods[55], a standard approach for interpreting the contribution of input features

643  to a neuron's activation (Figure 5F). Concretely, this technique produces an

644  attribution map for a given input cell, indicating each gene's contribution to that

645  cell's activation pattern. This framework enables single-cell or cell-state–level

646     gene explanations in a context-dependent manner. For instance, neuron 46

647     (N46) is activated in both epithelial cells and the ENS cells (Figure 5F), but its

648     activity is attributed to two distinct sets of highly co-expressed genes for each

649     context (Figure 5F, Supplemental Figure 7C). In the epithelial context, the top

650     genes for N46 activation are differentially expressed in goblet cells and enriched

651     for glycosylation pathways linked to goblet cell function. Meanwhile, in the ENS

652     context, the activating genes are highly co-expressed in late-stage neurons with

653     neuronal functions. Neither gene set shows strong expression outside its

654     respective context, indicating that CONCORD latent captures interpretable,

655     context-dependent gene co-expression programs.

656

**Discussion**

657

658     Contrastive learning has proven to be a powerful approach in image and

659     language processing, and recent work has begun to explore its adaptation for

660     single-cell sequencing data[17-23]. While some approaches use a supervised

661     paradigm[22,23], relying on annotated positive and negative samples to guide cell-

662     type separation, such strategies require substantial labeling and can be prone to

663     inaccuracies—particularly in single-cell contexts where cell identities may not be

664     completely defined. Other methods[17-20] adopt self-supervised contrastive learning

665     but employ a uniform sampling scheme that struggles to capture local

666     neighborhood structure and may inadvertently amplify batch effects. Additionally,

667     many of these methods use deep neural network architectures that require large

668     training datasets, limiting their applicability to single-cell studies with moderate

669     sample sizes.

670

671     In this work, we adapt self-supervised contrastive learning as a general

672     framework for dimensionality reduction and data integration in single-cell

673     analysis, operating in a fully unsupervised manner without requiring extensive

674     priors or large-scale training corpora. Central to our approach is a probabilistic

675     sampler that strengthens the contrastive learning process. Rather than drawing

676   cells entirely at random, our sampler selects from both global distributions and

677   local neighborhoods, enabling the model to observe broad variability across cell

678   types while also focusing on finer distinctions within closely related states. When

679   combined with masking and contrastive objectives, this design effectively learns

680   both large-scale and fine-grained gene co-expression patterns, filters out

681   inconsistent noise, and produces a high-resolution latent representation that

682   significantly improve downstream analyses, such as clustering, trajectory

683   inference, and UMAP visualizations.

684

685   A second key contribution of CONCORD is its capacity to integrate data across

686   multiple batches—an essential feature for large-scale efforts like the Human Cell

687   Atlas[56]. Conventional approaches often assume overlapping states or particular

688   batch-distortion models, leading to over- or under-correction that can disrupt

689   continuous trajectories or obscure biological signals[13]. In contrast, CONCORD

690   mitigates batch effects directly through its sampling strategy: the dataset-aware

691   sampler ensures that mini-batches primarily contain cells from a single dataset,

692   allowing the model to learn biologically meaningful variation while minimizing

693   dataset-specific artifacts. As a result, CONCORD generates a batch-effect-

694   mitigated latent representation even when overlap in cell states is minimal or

695   absent. Rather than stitching datasets together based on mutual-nearest-

696   neighbor-based 'anchors' or predefined reference points—a strategy prone to

697   accumulating distortions—CONCORD aligns cells based on underlying gene co-

698   expression structures, preserving biological continuity and structural integrity.

699   Moreover, its generalized sampling framework unifies neighborhood-aware and

700   dataset-aware sampling, enabling both effective data integration and enhanced

701   resolution. Crucially, these improvements demonstrate that contrastive learning

702   can be significantly enhanced without requiring complex model architectures or

703   additional training objectives, but rather through principled sampling and training.

704

Simulations and real-data benchmarks illustrate CONCORD's effectiveness in capturing fine-scale manifold structures. In simulated datasets, with or without batch effects, CONCORD consistently produces a denoised latent space that preserves critical topological and geometric characteristics. Importantly, these findings extend to real biological data. For instance, in embryonic atlases of *C. elegans* and *C. briggsae*, whose lineage trees are highly conserved, CONCORD not only integrated data from multiple species and batches spanning embryogenesis, but also accurately reconstructed intricate fate bifurcations and lineage convergences. Other methods either fell short in aligning these datasets, overlooked detailed structures, or fractured continuous trajectories. While narrowing the dataset to specific cell types and re-running integration can partially address these issues, it fragments any organism-wide perspective. CONCORD, by contrast, offers a unified, high-resolution atlas of developmental processes, enabling detailed tracing from progenitor cells to terminal states. Similar advantages emerged in mouse intestinal development, where CONCORD revealed complex hierarchies of subtypes, spatial zonation patterns, and cell cycle loops—all within a single integrated analysis. Unlike traditional workflows that remove cell cycle effects, CONCORD's embedding preserves both cell cycle loops and differentiation trajectories, allowing researchers to investigate the interplay between cell cycle and fate specification. In the intestine, analysis of the CONCORD latent revealed that differentiation trajectories are better described as complex topological structures where loops and cylinders are linked together by filaments. Furthermore, the simplicity of the CONCORD model enhances interpretability. Gradient-based attribution techniques allow users to identify the gene programs driving latent neuron activations, conferring mechanistic insights at single-cell or cell-type resolution.

The current CONCORD implementation employs a simple two-layer neural network, which significantly reduces memory requirements, improves computational efficiency, and scales to atlas-scale datasets while remaining robust for small datasets. However, its effectiveness may be limited in cases

736   where gene co-expression structures are highly distorted—such as in whole-cell

737   versus single-nucleus scRNA-seq—or when dropout rates vary significantly

738   across technologies. As with most neural network models, hyperparameter

739   selection is crucial. While CONCORD's minimalistic design reduces the number

740   of tunable parameters, key factors—including neighborhood size, neighborhood

741   enrichment probability, masking fraction, and contrastive temperature—can

742   substantially influence the final latent representation. Our benchmarking analysis,

743   along with prior studies[57], offers guidance on optimizing these parameters to

744   balance local resolution with global structural preservation and we offer a set of

745   default parameters that have performed well across multiple contexts.

746   Additionally, detailed tutorials on the CONCORD website

747   (https://qinzhu.github.io/Concord_documentation/) support users in making

748   informed hyperparameter choices. Future improvements to our sampling

749   strategy, particularly neighborhood-aware sampling, could incorporate adaptive

750   sampling probabilities that scale with distance, eliminating the need for

751   predefined neighborhood radii and enrichment probabilities. In its current form,

752   CONCORD supports a decoder for gene-level batch correction and a classifier

753   for tasks such as cell-type annotation or doublet detection. Early results suggest

754   these modules benefit from CONCORD's robust latent representation, though

755   additional optimization and benchmarking are required.

756

757   While the present benchmarks focused on single-cell RNA-seq, we have

758   observed promising outcomes in other single-cell modalities, including spatial

759   transcriptomics and scATAC-seq. Owing to its domain-agnostic design and

760   generalized sampler framework, CONCORD may extend beyond single-cell

761   transcriptional biology, offering a flexible, powerful approach for other single-cell

762   modalities and beyond.

763

764

765

**Methods**

**Self-supervised contrastive learning and sparse coding**

We implemented CONCORD in PyTorch, using a self-supervised contrastive learning approach inspired by SimCLR[14] and SimCSE[15], but with a unique dataset- and neighborhood-aware sampler design. The primary goal is to obtain a denoised, batch-effect-corrected latent representation that accurately captures essential characteristics of the underlying cell-state landscape. To achieve this, we apply the Normalized Temperature-scaled Cross Entropy Loss (NT-Xent loss)[14,15,58] directly to masking-augmented cell representations, following the design principles of unsupervised SimCSE[15].

A key feature of our approach is the integration of LeakyReLU activation functions with random masking augmentation, which enhances the capture of highly correlated gene co-expression patterns while suppressing spurious noise. Theoretical results[30] suggest that, if gene expression data $x$ can be approximated by a sparse coding model:

$$x = Mz + \varepsilon$$

where $Mz$ represents the sparse signal with $\|z\|_0 = \tilde{O}(1)$, and $\varepsilon$ denotes noise, then contrastive learning provably recovers the sparse features when trained with ReLU networks and random masking augmentation.

This sparse coding framework provides a generalized approach for modeling gene expression compared to widely used methods such as non-negative matrix factorization (NMF), principal component analysis (PCA), and factor analysis. Unlike these methods, sparse coding:

- Does not enforce orthogonality on $M$ (as in PCA),
- Does not require non-negativity constraints (as in NMF),
- Does not assume a probabilistic generative model (as in factor analysis).

At its core, this framework assumes an intrinsic low-rank latent structure, where gene co-expression programs define functional states—a phenomenon widely observed in single-cell transcriptomic studies[6,31,32]. Since contrastive learning has

27

795  been theoretically proven to capture such low-rank structures—or co-expressed

796  gene programs—we adopt a similar architecture in CONCORD and observe an

797  effective encoding of gene co-expression modules alongside noise suppression.

798  By relaxing constraints on orthogonality, non-negativity, and Gaussian noise

799  assumptions, the model may better capture diverse gene regulatory programs,

800  including those that do not conform to conventional constraints. Additionally,

801  random masking-based contrastive learning enhances robustness to dropout

802  artifacts and improves biological interpretability, ensuring that the learned latent

803  space faithfully represents gene programs underlying both discrete cell states

804  and continuous trajectories.

805

806  **Model architecture**

807  A schematic overview of the architecture is illustrated in Supplemental Figure 1

808  and described below:

809  1. Encoder:

810  The encoder takes randomly masked gene-expression vectors as input

811  and produces a latent encoding. By default, it is a two-layer fully

812  connected network (though the user can adjust the number of layers and

813  neurons). Optionally, a feature masking module with learnable weights can

814  be placed before the encoder to weight genes differently and encourage

815  sparse feature usage. Typically, single-cell data are normalized by total

816  count and log-transformed prior to training, though we have observed that

817  CONCORD remains robust to various normalization schemes as long as

818  no additional negative or zero values are introduced, and normalization is

819  applied consistently across datasets. During training, each cell is masked

820  at the gene level with a user-defined dropout probability (commonly 0.3–

821  0.6). Two masked versions of the same cell are fed through the encoder,

822  and their resulting embeddings are used in the contrastive loss.

823

824  2. Layer normalization and activation:

825    Each linear layer in the encoder is followed by layer normalization and a

826    user-selectable activation function (default: Leaky ReLU). Layer

827    normalization normalizes activations within each sample, making it more

828    robust to cross-batch variability than batch normalization[59] (which is also

829    implemented in CONCORD and available to user).

830

831    3.  Contrastive objective:

832    We employ the Normalized Temperature-scaled Cross Entropy Loss (NT-

833    Xent loss)[14,15,58] operating on mini-batches of N cells. Each cell is

834    randomly masked and encoded twice, yielding two different latent

835    representations: $z_i$ and $z_j$, then the contrastive loss is:

$$l_{i,j} = -\log\left(\frac{\exp\left(sim(z_i, z_i')/\tau\right)}{\sum_{k=1}^{2N} 1_{[k \neq i]} \exp\left(sim(z_i, z_k)/\tau\right)}\right)$$

837    where $sim(z_i, z_i') = \frac{z_i^T z_j}{\|z_i\|\|z_j\|}$ is the cosine similarity, and $\tau$ (default 0.5) is a

838    customizable temperature hyperparameter that controls the extent of local

839    separation and global uniformity of the embeddings[57].

840

841    4.  Decoder and classifier (optional):

842    For tasks requiring gene-expression reconstruction, an optional decoder

843    can be attached to the latent embeddings. In such cases, the latent vector

844    is combined with a learnable dataset embedding before being fed to the

845    decoder. This design separates batch information from the core latent

846    representation, preventing re-introduction of dataset-specific artifacts.

847

848    A classification head (i.e., a multi-layer perceptron) with cross-entropy loss

849    can be appended to the encoder for downstream tasks such as cell-type

850    prediction and doublet detection. The classifier can be trained on a pre-

851    trained encoder or jointly with the encoder to simultaneously guide cell-

852    type separations in the latent space. However, the latter approach may

853    impose a strong prior on the latent structure, potentially disrupting the

29

854      continuity of cellular trajectories. To mitigate overfitting in classification

855      tasks, a standard train-validation split with early stopping can be

856      employed.

857

858    **Dataset and neighborhood-aware probabilistic sampler**

859    The key innovation of CONCORD is the probabilistic mini-batch sampler, which

860    determines how cells are grouped and contrasted during training. Instead of

861    sampling cells uniformly at random, we introduce a generalizable sampler

862    framework that simultaneously (i) enriches local neighborhoods and (ii) focuses

863    each mini-batch primarily on a single dataset. This approach improves resolution

864    of detailed biological structures and yields a coherent, batch-effect-mitigated cell

865    atlas.

866

867    We begin by approximating the global data manifold with a k-nearest neighbors

868    (kNN) graph, where k is user-defined (and typically moderately large). This graph

869    can be built from normalized gene-expression data, a PCA projection, or a

870    preliminary CONCORD embedding. For large datasets, we employ the Faiss

871    library[60] for efficient neighbor searches. The kNN graph then facilitates sampling

872    at a user-specified local neighborhood enrichment probability, $p_{kNN}$ (default 0.3).

873

874    To construct mini-batches that are both dataset- and neighborhood-enriched, we

875    partition each mini-batch into four subsets—in-dataset neighbors, in-dataset

876    global samples, out-of-dataset neighbors, and out-of-dataset global samples

877    (Figure 1F). We randomly select a "core sample" from one dataset to anchor both

878    neighborhood selection and dataset enrichment. Then, each partition is sampled

879    based on $P_d$ (default 0.95) and $P_{kNN}$ as follows:

880    • In-dataset neighbors ($P_d \times P_{kNN}$): cells from the same dataset as the core

881      sample and lie in its KNN neighborhood.

882    • In-dataset global samples ($P_d \times (1 - P_{kNN})$): cells from the same dataset but

883      sampled uniformly from the entire data.

884   • Out-of-dataset neighbors ($(1 - P_d) \times P_{kNN}$):): cells from other datasets that

885       appear in the core sample's local kNN neighborhood.

886   • Out-of-dataset global samples ($(1 - P_d) \times (1 - P_{kNN})$): Cells from other

887       datasets, drawn uniformly across the entire data.

888

889   The sampler uses vectorized operations in PyTorch and NumPy for neighbor

890   retrieval, shuffling, and batch construction, thereby minimizing computational

891   overhead.

892

893   **Model training**

894   During each training epoch, CONCORD generates mini-batches using the

895   dataset and neighborhood-aware sampler, and randomly shuffles the mini-

896   batches. The Adam optimizer[61] was used to optimize the core contrastive

897   objective (NT-Xent), with optional loss terms—including reconstruction loss

898   (MSE) for the decoder, cross-entropy loss for classification, or feature-masking

899   penalties (L1/L2)—which can be incorporated as needed. By default, all loss

900   components are weighted equally but can be adjusted based on user-defined

901   preferences. A learning rate scheduler gradually reduces the learning rate to

902   enhance stability and convergence.

903

904   **Simulation pipeline**

905   We developed a versatile simulation pipeline to generate synthetic single-cell

906   gene expression data with diverse underlying structures. Unlike conventional

907   simulators that predominantly produce discrete clusters, our pipeline

908   accommodates a broad range of topologies, including linear trajectories,

909   branching trees, loops, and intersecting paths frequently observed in real single-

910   cell datasets.

911 In the first stage the state simulator constructs data according to a user-defined
912 structure:

913 • Clusters: Cells form discrete groups characterized by unique gene
914 programs, optionally including shared or ubiquitously expressed genes.

915 • Trajectories: Cells exhibit gradual shifts in gene expression, emulating cell
916 differentiation processes.

917 • Loops and intersecting paths: Continuous trajectories that close into loops
918 or intersect, representing cyclic biological processes.

919 • Trees: Hierarchical, branching lineages representing progenitor-to-terminal
920 fate differentiation, configurable by branching factor and tree depth.

921

922 With the chosen structure, the pipeline first generates a *noise-free* data matrix
923 with customizable cell and gene numbers. Expression values are then sampled
924 from selected distributions (e.g., Normal, Poisson, Negative Binomial),
925 introducing realistic variability and dropout patterns. Users can precisely control
926 parameters including mean baseline expression, dispersion (noise level), dropout
927 probability, and can enforce non-negativity or integer rounding of the generated
928 values.

929

930 In the second step, an optional batch simulator introduces dataset-specific
931 technical variability. This stage enables simulation of batch effects through
932 scaling factors, differential sampling rates, batch-specific gene subsets, and
933 expression-dependent dropout mechanisms. Multiple simulated batches are then
934 concatenated into a single dataset, with customizable proportions and varying
935 degrees of batch overlap to mimic real-world sampling scenarios.

936

937 By combining diverse gene expression structures with realistic noise models and
938 customizable batch effects, this simulation pipeline can approximate a broad

939 spectrum of biological and technical scenarios. As such, it provides a powerful

940 testbed for benchmarking data-integration techniques, trajectory-inference

941 algorithms, and manifold-learning methods under controlled yet biologically

942 realistic conditions.

943

944 **Benchmarking pipeline**

945 To comprehensively evaluate the performance of CONCORD and other

946 dimensionality reduction or data integration methods, we designed a robust

947 benchmarking pipeline that integrates geometric, topological, biological, and

948 batch-mixing metrics. This multifaceted assessment framework consists of the

949 following components:

950    1.  Topological assessments:

951       To quantify the preservation of intrinsic topological features, we employed

952       persistent homology analysis implemented via Giotto-TDA[62]. Persistent

953       homology captures structural properties of the data across multiple scales,

954       using Vietoris-Rips complexes constructed over increasing radii to

955       generate persistence diagrams and Betti curves. Persistence diagrams

956       reveal the lifespan of topological features such as connected components

957       (Betti-0), loops (Betti-1), and voids (Betti-2). We summarized these

958       diagrams through Betti curves and compared their mode—representing

959       the most persistent Betti number across all scales—to the known

960       topological ground truths. Additionally, we computed the entropy of Betti

961       curves to quantify the stability and complexity of the inferred topology, with

962       lower entropy reflecting stable and distinct structures, and higher entropy

963       indicating noisy or unstable topologies. These metrics were scaled

964       between 0 and 1 using min-max normalization to facilitate comparisons

965       among methods.

966

967    2.  Geometric assessments:

968    We evaluated the preservation of geometric relationships by calculating
969    distance correlations between embeddings and the corresponding noise-
970    free reference data, averaging Pearson, Spearman, and Kendall's tau
971    correlations to robustly quantify global geometric similarity. For local
972    neighborhood preservation, we employed trustworthiness[36], a metric
973    assessing how faithfully high-dimensional neighborhood structures are
974    maintained in lower-dimensional embeddings. Trustworthiness scores
975    range from 0 (poor preservation) to 1 (perfect preservation), and we
976    computed average trustworthiness scores across neighborhood sizes (k-
977    values) from 10 to 100 in increments of 10. Additionally, we visualized
978    trustworthiness as a function of k to reveal how each method performs at
979    different local scales. In cluster simulations with cluster-specific noise, we
980    further assessed the correlation of variance between the latent embedding
981    and the noisy input data, quantifying how accurately each method
982    preserves relative noise levels.

983

984    3. Batch mixing and biological label conservation:
985    We adopt established metrics from the scIB-metrics package[35] to
986    systematically evaluate biological label conservation and batch mixing.
987    Biological label conservation was quantified using metrics such as
988    Normalized Mutual Information (NMI) and Adjusted Rand Index (ARI),
989    which measure the correspondence between known biological labels and
990    inferred clusters. Batch mixing quality was assessed using silhouette
991    scores and graph metrics to determine how effectively embeddings
992    integrate cells from different batches or datasets.

993    A key limitation of the scIB pipeline is that the pipeline does not fully
994    accommodate the hierarchical and continuous nature of many biological
995    systems. Consequently, for simulations with continuous trajectories or
996    loops, we first apply Leiden clustering to noise-free data to define
997    "clusters" as ground truth, or use "branch" labels as a proxy for cell states

34

998     in tree simulations. Under these conditions, the scIB metrics are applied in

999     a more coarse-grained manner, offering an approximate assessment in

1000    these more complex scenarios.

1001

1002    **Data Availability**

1003    The human lung and pancreas datasets were compiled by Luecken et al.[35], and

1004    obtained from the scIB-metrics website (https://scib-

1005    metrics.readthedocs.io/en/stable/notebooks/lung_example.html) and the Open

1006    Problems in Single-Cell Analysis website

1007    (https://openproblems.bio/datasets/openproblems_v1/pancreas), respectively.

1008    The *C. elegans* embryogenesis atlas was downloaded from the Gene Expression

1009    Omnibus (GEO) (www.ncbi.nlm.nih.gov/geo) under accession code GSE126954.

1010    The joint *C. elegans* and *C. briggsae* dataset was obtained via email request

1011    from the authors of Large et al[41]. The mouse intestinal developmental atlas was

1012    acquired from GEO under accession code GSE233407.

1013

1014    **Code Availability**

1015    Concord is available at https://github.com/Gartner-Lab/Concord under a Creative

1016    Commons Attribution 4.0 International License. All benchmarking codes to

1017    generate results in this manuscript are deposited to https://github.com/Gartner-

1018    Lab/Concord_benchmark. Full documentation of Concord can be found at

1019    https://github.com/Gartner-Lab/Concord_documentation and available online at:

1020    https://qinzhu.github.io/Concord_documentation/.

1021

1022

## References

1. Wagner, D. E. & Klein, A. M. Lineage tracing meets single-cell omics: opportunities and challenges. *Nature Reviews Genetics* **21**, 410-427 (2020).

2. Tanay, A. & Regev, A. Scaling single-cell genomics from phenomenology to mechanism. *Nature* **541**, 331-338 (2017).

3. Waddington, C. H. *The strategy of the genes*. (Routledge, 2014).

4. Flores-Bautista, E. & Thomson, M. Unraveling cell differentiation mechanisms through topological exploration of single-cell developmental trajectories. *bioRxiv*, 2023.07. 2028.551057 (2023).

5. Johnson, J. A. *et al.* Inferring cellular and molecular processes in single-cell data with non-negative matrix factorization using Python, R and GenePattern Notebook implementations of CoGAPS. *Nature protocols* **18**, 3690-3731 (2023).

6. Kunes, R. Z., Walle, T., Land, M., Nawy, T. & Pe'er, D. Supervised discovery of interpretable gene programs from single-cell data. *Nature Biotechnology* **42**, 1084-1095 (2024).

7. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature methods* **16**, 1289-1296 (2019).

8. Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nature biotechnology* **37**, 685-691 (2019).

9. Stuart, T. *et al.* Comprehensive integration of single-cell data. *cell* **177**, 1888-1902. e1821 (2019).

10. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nature methods* **15**, 1053-1058 (2018).

11. Welch, J. D. *et al.* Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell* **177**, 1873-1887. e1817 (2019).

12. Haghverdi, L., Lun, A. T., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nature biotechnology* **36**, 421-427 (2018).

13. Zhang, Z. *et al.* Recovery of biological signals lost in single-cell batch integration with CellANOVA. *Nature Biotechnology*, 1-17 (2024).

14. Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. in *International conference on machine learning*. 1597-1607 (PMLR).

15. Gao, T., Yao, X. & Chen, D. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821* (2021).

16. He, K., Fan, H., Wu, Y., Xie, S. & Girshick, R. in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9729-9738.

17. Richter, T., Bahrami, M., Xia, Y., Fischer, D. S. & Theis, F. J. Delineating the effective use of self-supervised learning in single-cell genomics. *Nature Machine Intelligence*, 1-11 (2024).

18      Ciortan, M. & Defrance, M. Contrastive self-supervised clustering of scRNA-seq data. *BMC bioinformatics* **22**, 280 (2021).

19      Wang, J., Xia, J., Wang, H., Su, Y. & Zheng, C.-H. scDCCA: deep contrastive clustering for single-cell RNA-seq data based on auto-encoder network. *Briefings in Bioinformatics* **24**, bbac625 (2023).

20      Zhao, B., Song, K., Wei, D.-Q., Xiong, Y. & Ding, J. scCobra allows contrastive cell embedding learning with domain adaptation for single cell data integration and harmonization. *Communications Biology* **8**, 233 (2025).

21      Yang, M. *et al.* Contrastive learning enables rapid mapping to multimodal single-cell atlas of multimillion scale. *Nature Machine Intelligence* **4**, 696-709 (2022).

22      Heimberg, G. *et al.* A cell atlas foundation model for scalable search of similar human cells. *Nature*, 1-3 (2024).

23      Heryanto, Y. D., Zhang, Y.-z. & Imoto, S. Predicting cell types with supervised contrastive learning on cells and their types. *Scientific Reports* **14**, 430 (2024).

24      Yang, Z. *et al.* in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.  3057-3069.

25      Goodfellow, I. *et al.* Generative adversarial networks. *Communications of the ACM* **63**, 139-144 (2020).

26      Lotfollahi, M., Wolf, F. A. & Theis, F. J. scGen predicts single-cell perturbation responses. *Nature methods* **16**, 715-721 (2019).

27      Ganin, Y. & Lempitsky, V. in *International conference on machine learning*. 1180-1189 (PMLR).

28      Sohn, K., Lee, H. & Yan, X. Learning structured output representation using deep conditional generative models. *Advances in neural information processing systems* **28** (2015).

29      Riba, A. *et al.* Cell cycle gene regulation dynamics revealed by RNA velocity and deep-learning. *Nature communications* **13**, 2865 (2022).

30      Wen, Z. & Li, Y. in *International Conference on Machine Learning*.  11112-11122 (PMLR).

31      Jiang, J. *et al.* D-SPIN constructs gene regulatory network models from multiplexed scRNA-seq data revealing organizing principles of cellular perturbation response. *BioRxiv*, 2023.2004. 2019.537364 (2024).

32      Kotliar, D. *et al.* Identifying gene expression programs of cell-type identity and cellular activity with single-cell RNA-Seq. *Elife* **8**, e43803 (2019).

33      Alaqeeli, O. A comparison of dropout rate of three commonly used single cell RNA-sequencing protocols. *Biotechnology & Biotechnological Equipment* **38**, 2379837 (2024).

34      Zappia, L., Phipson, B. & Oshlack, A. Splatter: simulation of single-cell RNA sequencing data. *Genome biology* **18**, 174 (2017).

35      Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics. *Nature methods* **19**, 41-50 (2022).

1108    36    Venna, J. & Kaski, S. in *International conference on artificial neural networks*.
1109          485-491 (Springer).
1110    37    McInnes, L., Healy, J. & Melville, J. UMAP: uniform manifold approximation
1111          and projection for dimension reduction. arXiv. *arXiv preprint*
1112          *arXiv:1802.03426* **10** (2018).
1113    38    Van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *Journal of*
1114          *machine learning research* **9** (2008).
1115    39    Tenenbaum, J. B., Silva, V. d. & Langford, J. C. A global geometric framework
1116          for nonlinear dimensionality reduction. *science* **290**, 2319-2323 (2000).
1117    40    Sulston, J. E., Schierenberg, E., White, J. G. & Thomson, J. N. The embryonic
1118          cell lineage of the nematode Caenorhabditis elegans. *Developmental biology*
1119          **100**, 64-119 (1983).
1120    41    Large, C. R. *et al.* Lineage-resolved analysis of embryonic gene expression
1121          evolution in C. elegans and C. briggsae. *bioRxiv*, 2024.2002. 2003.578695
1122          (2024).
1123    42    Packer, J. S. *et al.* A lineage-resolved molecular atlas of C. elegans
1124          embryogenesis at single-cell resolution. *Science* **365**, eaax1971 (2019).
1125    43    Büttner, M., Miao, Z., Wolf, F. A., Teichmann, S. A. & Theis, F. J. A test metric
1126          for assessing single-cell RNA-seq batch correction. *Nature methods* **16**, 43-
1127          49 (2019).
1128    44    Tintori, S. C., Nishimura, E. O., Golden, P., Lieb, J. D. & Goldstein, B. A
1129          transcriptional lineage of the early C. elegans embryo. *Developmental cell*
1130          **38**, 430-444 (2016).
1131    45    Ortiz, C. O. *et al.* Searching for neuronal left/right asymmetry: genomewide
1132          analysis of nematode receptor-type guanylyl cyclases. *Genetics* **173**, 131-
1133          149 (2006).
1134    46    Yu, S., Avery, L., Baude, E. & Garbers, D. L. Guanylyl cyclase expression in
1135          specific sensory neurons: a new family of chemosensory receptors.
1136          *Proceedings of the National Academy of Sciences* **94**, 3384-3387 (1997).
1137    47    Koreth, J. & van den Heuvel, S. Cell-cycle control in Caenorhabditis elegans:
1138          how the worm moves from G1 to S. *Oncogene* **24**, 2756-2764 (2005).
1139    48    Huycke, T. R. *et al.* Patterning and folding of intestinal villi by active
1140          mesenchymal dewetting. *Cell* **187**, 3072-3089. e3020 (2024).
1141    49    Zwick, R. K. *et al.* Epithelial zonation along the mouse and human small
1142          intestine defines five discrete metabolic domains. *Nature Cell Biology* **26**,
1143          250-262 (2024).
1144    50    Morarach, K. *et al.* Diversification of molecularly defined myenteric neuron
1145          classes revealed by single-cell RNA sequencing. *Nature neuroscience* **24**,
1146          34-46 (2021).
1147    51    Derecka, M. *et al.*    (American Society of Hematology Washington, DC,
1148          2017).
1149    52    Chen, C.-P., Wang, L.-K., Chen, C.-Y., Chen, C.-Y. & Wu, Y.-H. Placental
1150          multipotent mesenchymal stromal cell-derived Slit2 may regulate

1151     macrophage motility during placental infection. *Molecular Human*
1152     *Reproduction* **27**, gaaa076 (2021).
1153  53 Al-Shboul, O. A. The importance of interstitial cells of cajal in the
1154     gastrointestinal tract. *Saudi Journal of Gastroenterology* **19**, 3-15 (2013).
1155  54 Torihashi, S. *et al.* Blockade of kit signaling induces transdifferentiation of
1156     interstitial cells of cajal to a smooth muscle phenotype. *Gastroenterology*
1157     **117**, 140-148 (1999).
1158  55 Ancona, M., Ceolini, E., Öztireli, C. & Gross, M. Gradient-based attribution
1159     methods. *Explainable AI: Interpreting, explaining and visualizing deep*
1160     *learning*, 169-191 (2019).
1161  56 Atlas, H. C. Human cell atlas. *Human cell atlas sequences first 250K*
1162     *developmental cells* (2018).
1163  57 Wang, F. & Liu, H. in *Proceedings of the IEEE/CVF conference on computer*
1164     *vision and pattern recognition.*  2495-2504.
1165  58 Sohn, K. Improved deep metric learning with multi-class n-pair loss
1166     objective. *Advances in neural information processing systems* **29** (2016).
1167  59 Lei Ba, J., Kiros, J. R. & Hinton, G. E. Layer normalization. *ArXiv e-prints*, arXiv:
1168     1607.06450 (2016).
1169  60 Douze, M. *et al.* The faiss library. *arXiv preprint arXiv:2401.08281* (2024).
1170  61 Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv*
1171     *preprint arXiv:1412.6980* (2014).
1172  62 Tauzin, G. *et al.* giotto-tda:: A topological data analysis toolkit for machine
1173     learning and data exploration. *Journal of Machine Learning Research* **22**, 1-6
1174     (2021).

1175

**Acknowledgement**

**Author Contributions**

Q.Z. and Z.J.G. conceived the project. Q.Z. designed and implemented the method. Q.Z. conducted benchmarking analyses with assistance from Z.J. M.T. provided critical feedback and methods for topological data analysis. Z.J.G. supervised the project. Q.Z. and Z.J.G. wrote the manuscript. All authors reviewed and edited the manuscript.

**Competing Interests**

ZJG is an author on patents associated with sample multiplexing and ZJG is an equity holder and advisor to Provenance Bio.

**Materials & Correspondence**

1203    Correspondence to Q. Zhu (qin.zhu@ucsf.edu) and Z. J. Gartner

1204    (zev.gartner@ucsf.edu).

1205

1206 **Figures**



**Figure 1. The CONCORD sampler enables contrastive learning to generate a high-resolution, batch-effect-mitigated latent representation of scRNA-seq data.**

**(A)** Illustration of hypothetical cell state landscapes and corresponding dimensionality-reduced representations that capture key structural features of the landscape. (**B**) Comparison of neighborhood-aware and uniform sampling and their impact on contrastive learning in a simulated four-cell-state dataset. The heatmap shows the actual simulated expression. For each sampling scheme, PCA plots are color-coded by cell state, with black points indicating cells selected in a representative mini-batch, accompanied by density curves illustrating their distribution. Latent heatmaps display the representations learned using uniform and neighborhood-aware sampling, with black lines marking cells included in the selected mini-batch and density plots depicting their distribution. The resulting UMAP embeddings computed from the latent representations for each sampling method are also shown. (**C**) Contrastive learning performed in a single batch with the conventional sampler, which draws cells uniformly from the entire dataset to form mini-batches. (**D**) When applying standard contrastive learning to multiple datasets (represented by the blue or pink background), contrasting cells from different datasets within the same mini-batch amplifies dataset-specific biases, which is manifested in the latent embeddings. (**E**) CONCORD mitigates these dataset-specific artifacts by predominantly contrasting cells within each dataset and randomly shuffling mini-batches for each training epoch. (**F**) The CONCORD sampling framework. A "leaky" dataset-aware sampler addresses minimal or absent overlap between datasets and can be combined with the neighborhood-aware sampler to support both data integration and enhanced resolution. This is achieved by a joint probabilistic sampling framework, where the likelihood of selecting a given cell reflects the combined probabilities of dataset-aware ($P_d$) and neighborhood-aware sampling ($P_{kNN}$).

1207

1208

1209

1210

1211

1212

**Figure 2. Benchmarking CONCORD and other dimensionality reduction methods across diverse structures.**

**(A)** Simulation pipeline for generating data structures. The pipeline first produces a noise-free gene expression matrix based on a user-defined structure, then introduces noise following a specified noise model, and finally

applies batch effects in various forms. (**B**) Evaluation pipeline. Using the simulated datasets, the latent representations produced by each method was compared with the noise-free ground truth to assess how well topological and geometric features are preserved. For cluster simulations, we further evaluate the correlation of cluster-specific variances in the noisy data versus the latent space. Metrics from the scIB[35] package were incorporated for evaluating conservation of biological labels and harmonization of batch effects. (**C**) Performance on simulated clusters, highlighting the resulting UMAP visualization, cosine distance matrices, persistence diagrams, and Betti curves for CONCORD and other methods. In the persistent homology analysis, the $H_0$ point representing infinity was excluded from the persistence diagram and curve. (**D**) Performance on a complex trajectory with 3 loops, highlighting the same diagnostic plots as C. (**E**) Summary table for the three-cluster simulation, listing key topological and geometric evaluation metrics. (**F**) Table summarizing the methods' performance on the complex trajectory-loop simulation. (**G**) KNN graph visualization of latent embeddings from each method on a complex tree simulation, with zoomed-in views of the darkened region highlighting detail on one of the branches.

1213

1214

**Figure 3. Benchmarking CONCORD and other data integration methods across diverse structures.**

(**A**) Two-batch, five-cluster simulation. For ground truth, we show kNN graphs (k=15 by default, with edges omitted) of both the noise-free and noise-added data (no batch effect). Latent spaces from each integration method are

visualized by kNN graphs, colored by batch (top) and cluster (bottom). (**B**) Trajectory simulation with varying batch overlap. The ground truth is shown with PCA and a kNN graph. For each method, the resulting latent space is depicted with a kNN graph (k = 15) to assess how well cells are integrated across batches along the trajectory. In the gap simulation, an additional kNN graph (k = 30), colored by simulated time, demonstrates that CONCORD accurately captures the correct orientation of the trajectories along time despite the gap. (**C**) Loop simulation with varying batch overlap. Shown here are kNN graphs of the ground truth (with edges omitted) and the CONCORD latent space. Full results for other methods are provided in Supplementary Figure 3D. (**D**) Tree simulation with varying batch overlap. kNN graphs for the ground truth and the CONCORD latent space are shown. Full results for other methods are provided in Supplementary Figure 3E. (**E**) scIB benchmarking on the two-batch, five-cluster simulation. Integration performance was evaluated using metrics from the scIB-metrics package[35]. (**F**) Ranking of integration methods. Each method's performance is scored across topological, geometric, and scIB metrics. The overall rank is based on the average ranking across all metrics.

1215

1216

**Figure 4. Benchmarking CONCORD on *C. elegans/C.briggsae* embryogenesis atlas.**

(**A**) Global 2D and 3D UMAPs of CONCORD (with decoder) colored by cell type and estimated embryo time. (**B**) UMAP of CONCORD and other integration methods colored by estimated embryo time and species. (**C**) Boxplots show the fraction of *C. elegans* cells within randomly sampled 100-

nearest-neighbor (100-NN) neighborhoods, stratified by embryo time bins. The red horizontal line represents the expected species fraction based on the global composition of each time bin. Well-integrated datasets should show species fractions closely matching this expected value with minimal variation. (**D**) Global 3D UMAPs of CONCORD (with decoder), Seurat and scVI, highlighting cells mapped to the lineage sub-tree that give rise to ASE, ASJ and AUA neurons. (**E**) Heatmap showing the top 50 most variable latent dimensions in the ASE, ASJ, and AUA neuron subset for scVI, Seurat, and CONCORD (with decoder). Expression of gcy-5 and gcy-14 were plotted on the CONCORD (with decoder) UMAP. (**F**) Latent space distance between medoids of ectodermal cells (AB lineage), stratified by cell generation and lineage relationship. AB5 refers to cells derived after five successive divisions of the AB founder cell, with AB6 to AB9 representing progressively later generations. (**G**) Spearman correlation between lineage distance and latent space distance across integration methods for AB lineages from generations 5 to 9. Statistical significance of differences in correlations was assessed using a two-sided Mann-Whitney U test, with asterisks indicating significance levels (\*\*$p < 0.01$, \*\*\*$p < 0.001$, \*\*\*\*$p < 0.0001$). (**H**) Zoom-in UMAPs for mesoderm cells excluding pharynx. Major input lineages and cell types were highlighted. Each lineage was represented by its cluster medoid on the UMAP, and lines connect each parental lineage to its daughter lineages following the lineage tree. (**I**) Zoom-in UMAPs for pharynx, annotated with cell types and broad input lineages. Selected lineage paths that give rise to pm1/2, pm3-5, and pm6 are highlighted. (**J**) Run time comparison of different integration methods. \*Harmony was run using a 300-dimensional PCA input, whereas all other methods were applied to the gene expression matrix containing 10,000 variably expressed genes.

1217

1218

49

**Figure 5. Benchmarking CONCORD on mammalian intestine development.**

(**A**) 2D and 3D UMAP visualizations of CONCORD latent space, colored by cell type and cell cycle phase, with cell-type-colored UMAPs from scVI and Seurat shown for comparison. (**B**) Zoom-in views of epithelial cells in the 3D global

UMAP, colored by cell subtype, zonation, and expression of zonation-specific markers (Bex1, Onecut2). A red marker and arrow indicate the viewing angle within the 3D global UMAP. (**C**) Zoom-in view of enteric nervous system (ENS) cells, colored by cell cycle phase and cell state/branch annotations, based on Morarach et al[50], along with state-specific gene expression. A red marker and arrow indicate the viewing angle. (**D**) Zoom-in view of Pdgfra- mesenchymal cells and smooth muscle cells, colored by cell cycle phase, subtype annotation, and selected subtype-specific markers. A red marker and arrow indicate the viewing angle. (**E**) Heatmap of all latent encodings generated by CONCORD, Seurat, and scVI. (**F**) Interpretation of CONCORD latent space using gradient-based attribution techniques. Activation of Neuron 46 (Z46) in epithelial and ENS cells is attributed to the co-expression of epithelial- and neuron-specific gene sets in their respective contexts. GO enrichment analysis of these gene sets is shown.

1219

1220 **Supplementary Information**

1221 **Supplementary Figures**



**Supplemental Figure 1. The CONCORD framework.**

Each cell undergoes random masking twice, and contrastive loss is computed on the latent encodings of mini-batches sampled using the dataset- and neighborhood-aware sampler. An optional decoder, incorporating a learnable dataset covariate, enables batch-free gene expression inference, while an optional classifier can be included for cell type classification or annotation-guided learning.

1222

1223

**Supplemental Figure 2. Benchmarking CONCORD and other dimensionality reduction methods across diverse structures.**

(**A**) Heatmaps of the simulated expression for the 3-cluster structure and the corresponding CONCORD latent encoding. (**B**) Heatmaps of the simulated expression for the trajectory-loop structure and the corresponding CONCORD latent encoding. (**C**) Trustworthiness measured across varying neighborhood sizes in the 3-cluster simulation. Note that in the noise-free setting, within-cluster neighbors are assigned randomly, so trustworthiness does not equal 1 for the no-noise data. (**D**) Trustworthiness measured across varying neighborhood sizes in the complex trajectory-loop simulation. (**E**) Heatmaps of the simulated expression for the complex tree structure shown in Figure 2G and the CONCORD latent encoding with different intra-kNN enrichments. (**F**)

Trustworthiness measured across varying neighborhood sizes in the complex tree simulation, evaluated under different intra-kNN enrichment conditions. A zoomed-in view of the region where k < 20 and trustworthiness > 0.96 highlights improvements in capturing local neighborhood relationships with moderate intra-kNN enrichments.
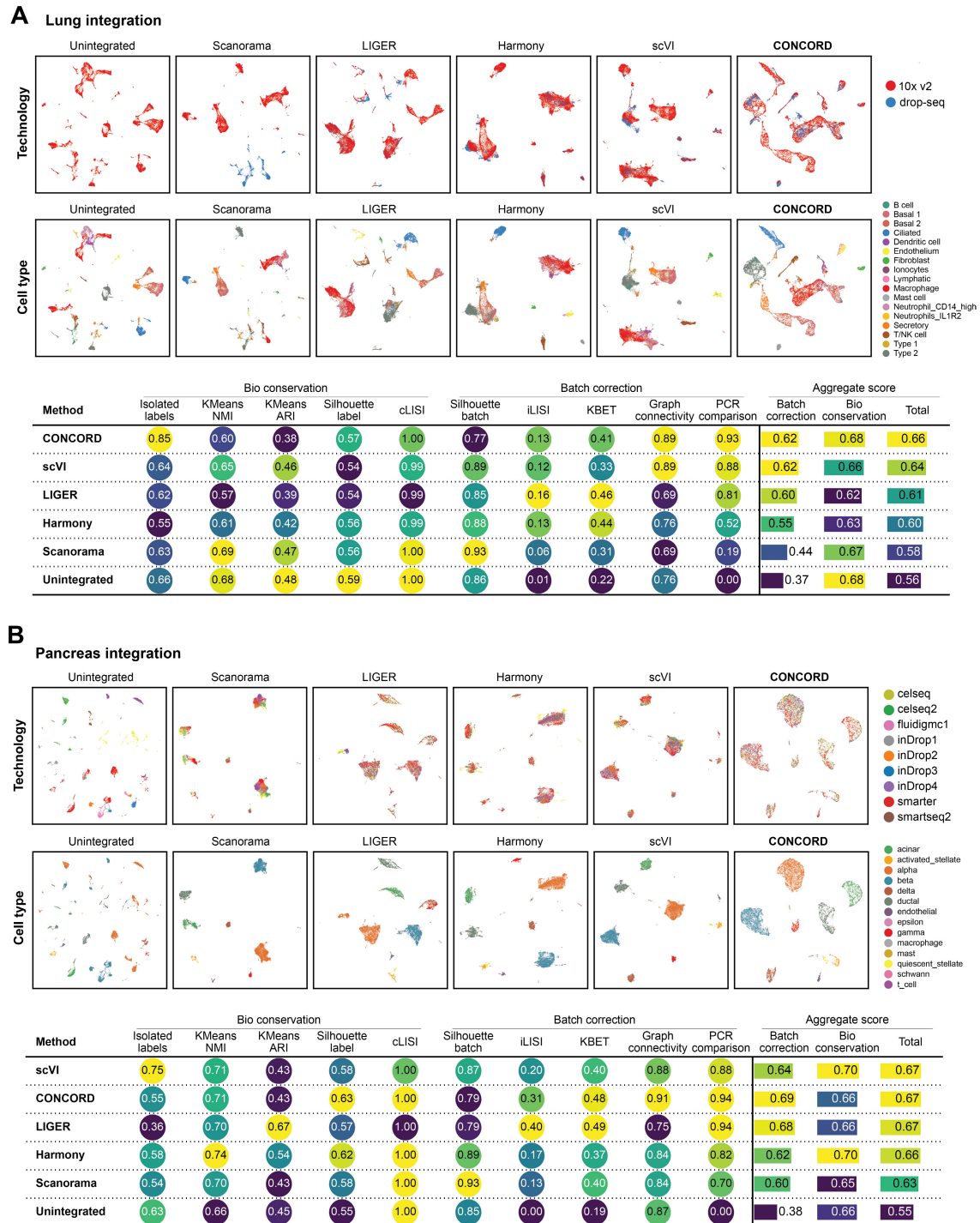
1224

1225

**Supplemental Figure 3. Benchmarking CONCORD and other data integration methods across diverse structures.** (**A**) Correlation of cluster-specific variances between noise-added ground truth and latent from each integration methods. (**B**) Two-batch, five-cluster simulation with minimal overlap. The two batches share only one cluster. kNN graphs (edges omitted) from each integration method are shown. (**C**) Trajectory simulation with full and

partial overlap. kNN graphs (edges omitted) of the ground truth and CONCORD latent with varying intra-dataset enrichment ($P_d$) are shown, colored by batch (top) and time (bottom). (**D**) Loop simulation with varying batch overlap. kNN graphs (edges omitted) 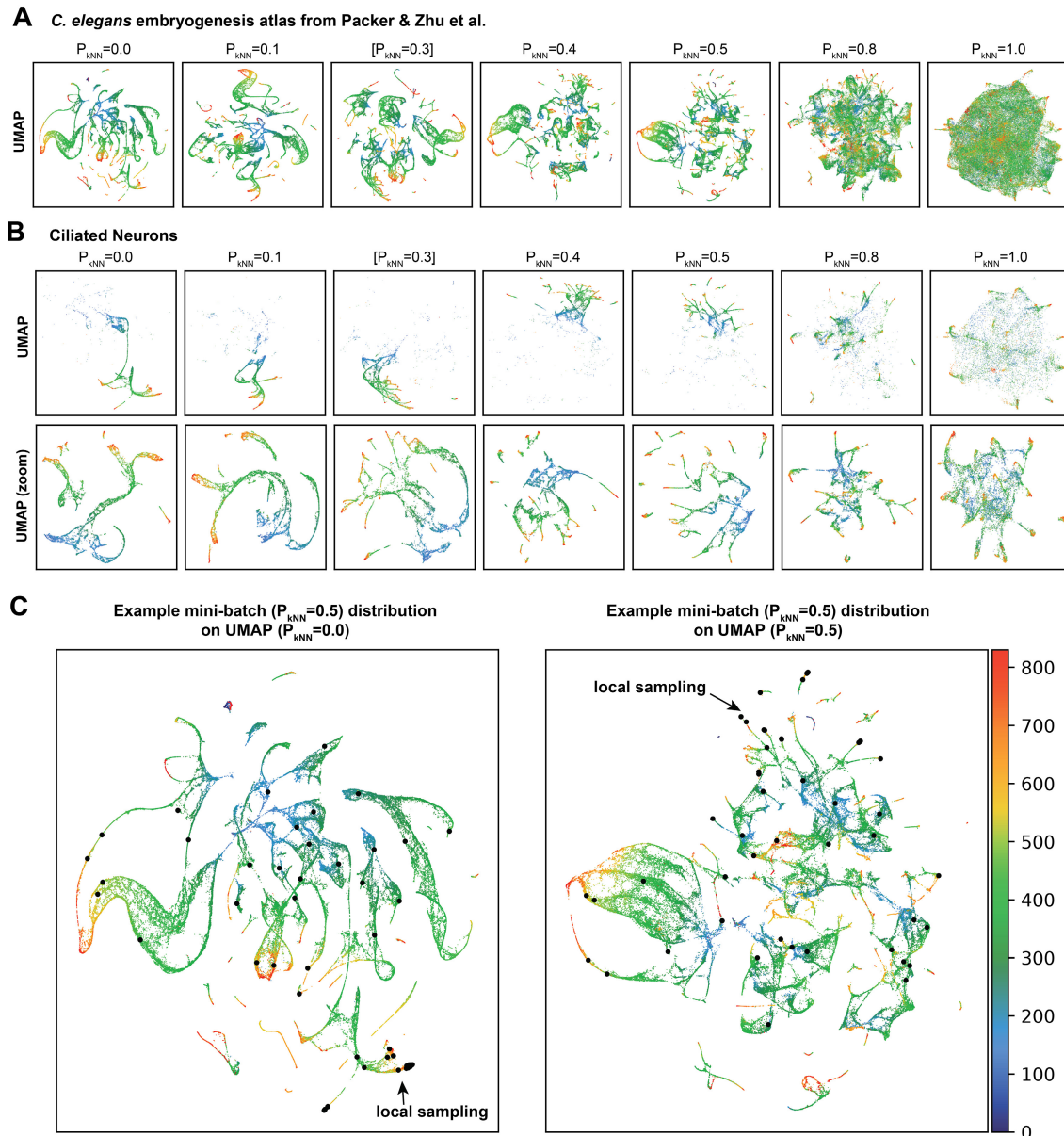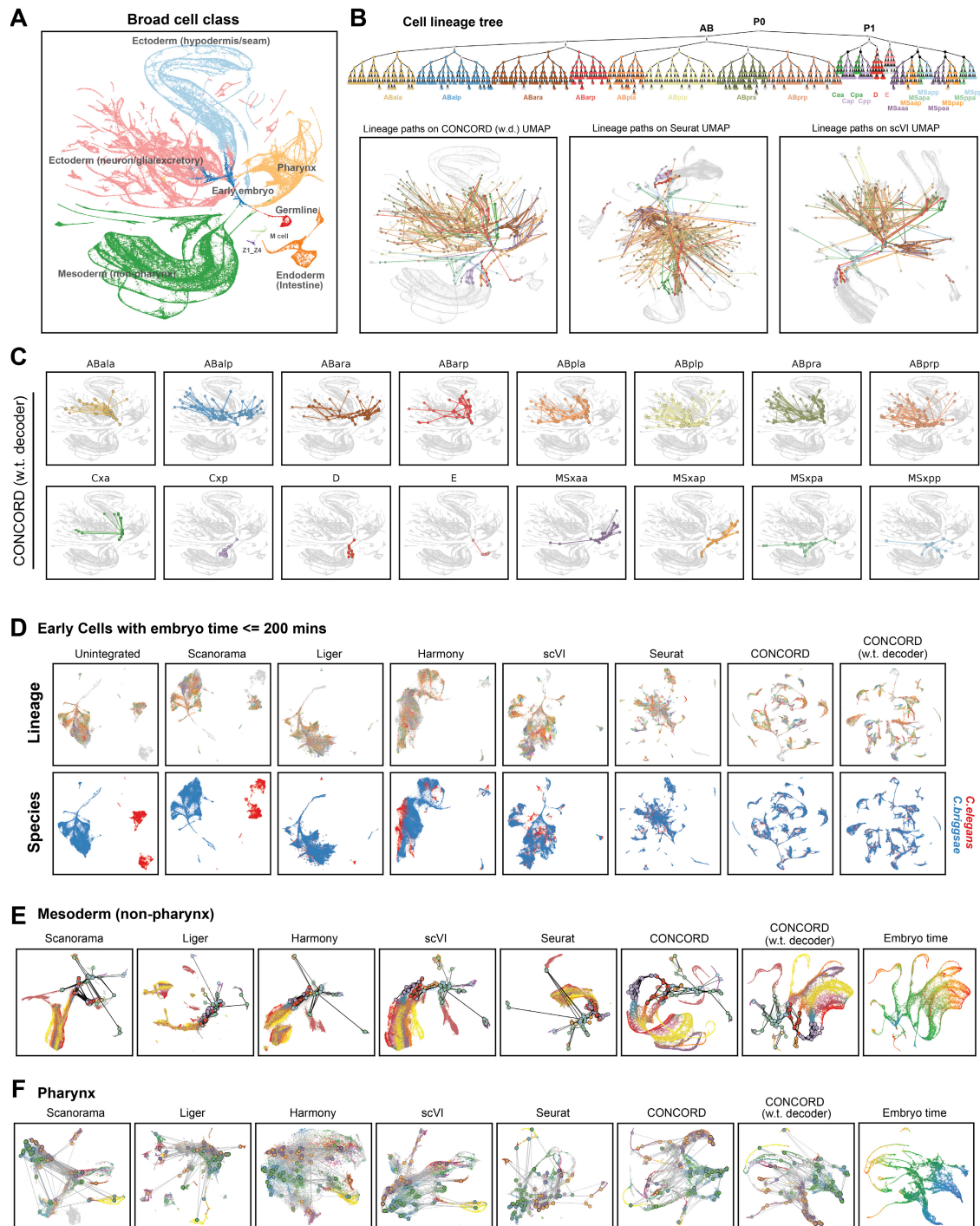of the ground truth and each integration method are shown. (**E**) Tree simulation with varying batch overlap. kNN graphs (edges omitted) of the ground truth and each integration method are shown.

1226

1227

1228

**Supplemental Figure 4. Performance of CONCORD on lung and pancreas datasets.** (**A**) UMAPs were generated for both unintegrated lung data[35] and the latent space from each integration method, colored by technology and cell type. Benchmarking statistics from the scIB-metrics package[35] were shown. (**B**) Benchmarking of CONCORD on the integration of pancreas dataset[35]. UMAPs were generated for both unintegrated data and the latent space from each

integration method, colored by technology and cell type. Benchmarking statistics from the scIB-metrics package[35] were shown.
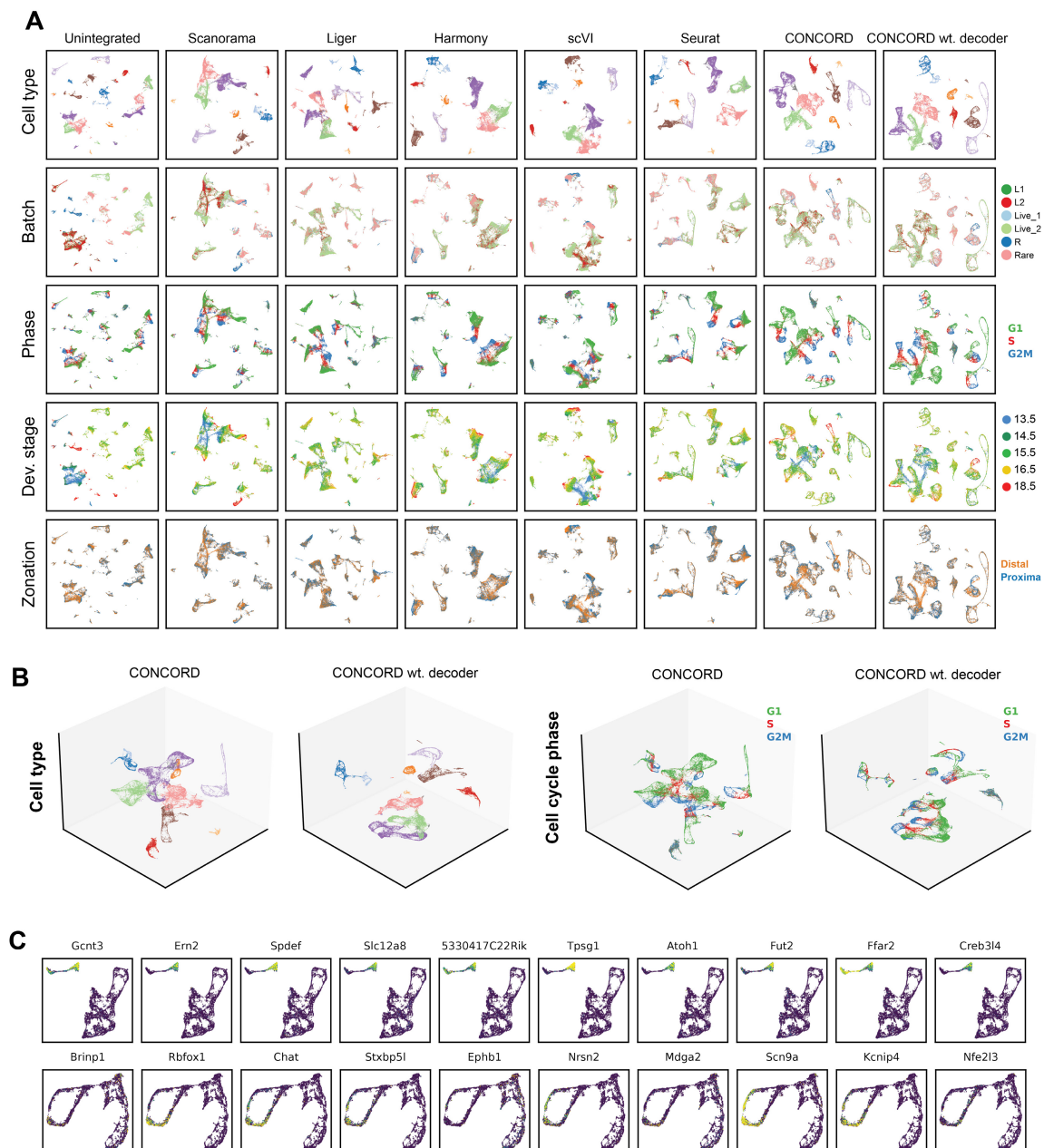
1229

**Supplemental Figure 5. Effect of neighborhood-aware sampling on *C. elegans* embryogenesis atlas.** (**A**) Global UMAP of CONCORD latent space with varying levels of intra-kNN enrichment, colored by embryonic time. (**B**) Global and zoomed-in UMAPs highlighting ciliated neurons, colored by embryonic time, generated from CONCORD latent representations across different levels of intra-kNN enrichment. (**C**) Visualization of cells sampled within a single mini-batch using an intra-kNN probability of 0.5. These cells are highlighted in the global UMAPs of CONCORD without neighborhood enrichment and with 0.5 neighborhood enrichment, demonstrating the effect of local sampling on structural resolution.
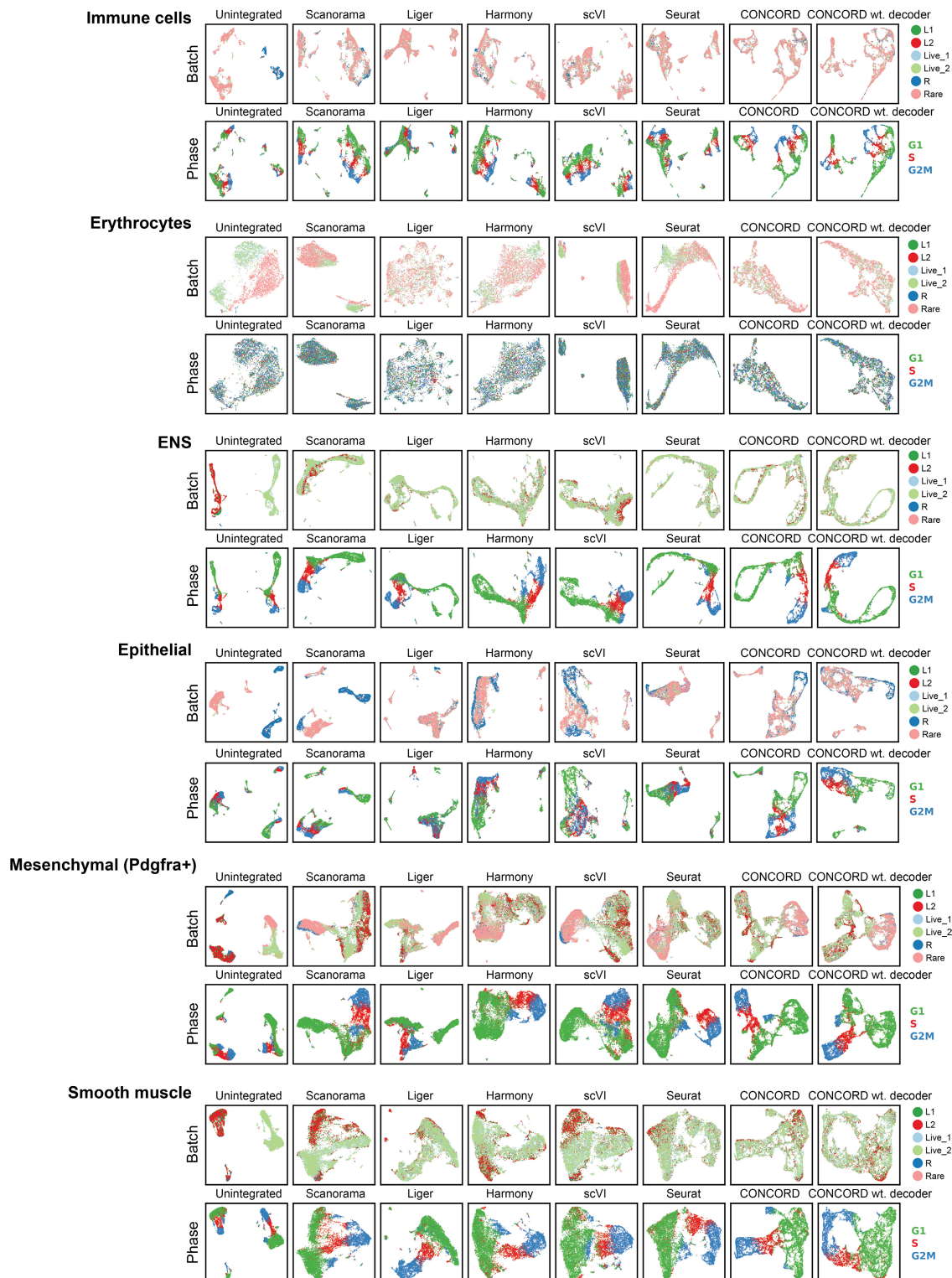
1230

**Supplemental Figure 6. Cell lineage tree and zoom-in analysis of C. elegans/C.briggsae embryogenesis atlas.** (**A**) Global UMAP of CONCORD (with decoder) colored by broad cell class. (**B**) The lineage annotations by Large et al. were mapped to the *C. elegans* lineage tree, with some ambiguous mappings due to symmetry. Each lineage was represented by its cluster medoid on the UMAP, and lines connect each parental lineage to its daughter lineages following the lineage tree. (**C**) For CONCORD (with decoder), lineage

sub-trees of major lineage groups were separately highlighted. (**D**) UMAPs of early embryo cells with embryo time before or equal to 200 minutes, colored by lineage and species. (**E**) Zoom-in UMAPs for mesoderm cells excluding pharynx. Each input lineage was represented by its cluster medoid on the UMAP, and lines connect each parental lineage to its daughter lineages following the lineage tree. Estimated embryo time was also plotted on the UMAP generated with CONCORD (with decoder) latent. (**F**) Zoom-in UMAPs for pharynx, colored by cell types and input lineages. Lineage paths were plotted following the lineage tree. Estimated embryo time was also plotted on the UMAP generated with CONCORD (with decoder) latent.

1231

**Supplemental Figure 7. Benchmarking CONCORD on mammalian intestine development.** (**A**) UMAPs of the mouse intestinal developmental atlas[48] generated from CONCORD latent space, other integration methods, and unintegrated data, colored by broad cell type, batch, cell cycle phase, developmental stage, and zonation. (**B**) 3D UMAPs of CONCORD with and without the decoder, colored by cell type and cell cycle phase. (**C**) Expression patterns of the top genes contributing to Neuron 46 activation in the epithelial context (top) and ENS context (bottom).

1232

**Supplemental Figure 8. Zoom-in UMAPs of major cell types in the mammalian intestinal developmental atlas.** Zoom-in UMAPs for each major cell type computed from CONCORD latent space, other integration methods, and unintegrated data, colored by batch (top) and cell cycle phase (bottom).

## Supplementary Tables

**Supplemental Table 1. Benchmarking data integration methods across diverse structures with varying degrees of overlap.** For each structure type (cluster, trajectory, loop, and tree), we simulated varying degrees of overlap (full overlap, partial overlap, connected without overlap, and gap between batches). The table displays the score of each method using topological, geometric, and scIB metrics[35].