# Teaching AI Agents to See: Geometric Machine Learning for Single-Cell Data

## A Machine Learning Approach to Algorithm Selection

**Brian Di Bassinga**

A senior thesis presented to the faculty of the
Department of Computer Science
Yale University
in partial fulfillment of the requirements for the
Bachelor of Science degree

Advisor: Professor Smita Krishnaswamy
Department of Computer Science
Yale University

December 2025

# Acknowledgements

# Contents

# Abstract

Multi-agent systems for biological data analysis increasingly rely on large language models (LLM) for coordination. However, the use of natural language communication by these LLMs introduces systematic information loss when agents must describe high-dimensional geometric properties. This thesis investigates whether geometric metrics computed from single-cell embeddings can predict structural categories—a prerequisite for geometry-informed algorithm selection in automated analysis pipelines.

We developed complete infrastructure for this investigation: a literature mining pipeline that extracted algorithm choices from 167 published papers, a geometric feature extraction system computing 53 structural metrics on embeddings, and a 7-class taxonomy for characterizing embedding geometry (clusters, multi-branch, horseshoe, bifurcation, diffuse, simple trajectory, cyclic). We processed 100 single-cell datasets from the CELLxGENE Census and created 100 expert annotations, resulting in 91 high-quality training samples after quality control.

Our classifier achieved 40.43% accuracy on the 7-class task ($2.9\times$ random baseline) using Support Vector Machine with RBF kernel. Feature importance analysis revealed that density-based metrics (spatial entropy, density coefficient of variation) and hull compactness were most predictive of embedding structure. The classifier successfully distinguished distinctive shapes like horseshoes and clusters but struggled with visually similar categories, such as clusters versus multi-branch structures.

These results demonstrate that geometric metrics correlate with embedding structure, establishing feasibility for geometry-informed multi-agent coordination. The infrastructure developed—including the ManyLatents framework for metric computation and the CELLxGENE integration pipeline—provides a foundation for scaling this approach with additional labeled data.

# 1   Introduction

## 1.1   Problem Statement

The rapid availability of high-dimensional biological data has created an urgent need for automated analytical frameworks. While systems such as CellAgent and Biomni have demonstrated proficiency in orchestrating general research tasks, they face a critical limitation when it comes to understanding the shapes and structures that emerge when high-dimensional genomic data is reduced to visual embeddings.

Current multi-agent systems rely exclusively on natural language to describe data structure. However, the precise geometric properties of an embedding—how dense different regions are, whether the data forms loops or branches, how it curves through space—cannot be accurately conveyed through words without significant information loss. A descriptor like 'branching trajectory' fails to distinguish between a single fork and a complex tree of branches, yet this distinction is crucial for selecting the appropriate dimensionality reduction algorithm.

Consequently, when agents coordinate complex analysis pipelines, they cannot reliably communicate what they've observed about data structure. One agent might describe an embedding as having 'distinct clusters,' but without quantitative details—how many clusters, how separated, how dense—a downstream agent cannot make informed decisions about which algorithms to apply next. This thesis argues that effective automated analysis requires agents to share precise, numeric descriptions of embedding structure rather than relying on vague linguistic approximations.

## 1.2   Research Questions

This thesis addresses three core questions:

1. **Can geometric metrics characterize embedding structure?** Do the quantitative properties of biological data (local intrinsic dimensionality, participation ratio, diffusion curvature, topological features) correlate with visually distinguishable dimensionality reduction embedding categories?

2. **Can machine learning automate geometry-based structure classification?** Can we train a classifier to predict embedding structure categories from computed geometric metrics?

3. **What infrastructure is needed for systematic investigation?** What tools and datasets are required to pursue geometry-informed algorithm selection systematically?

## 1.3   Significance and Impact

This research establishes the feasibility of geometric-feature-based structure classification and provides the infrastructure needed for geometry-informed analysis systems. The

contributions include:

**Methodological Infrastructure:** Developed a complete, scalable pipeline for geometry-informed analysis, including a literature mining system that extracted algorithm choices from 167 papers. Leveraging the lab's ManyLatents framework for computing 53 structural metrics, this infrastructure addresses the data engineering requirements for systematic geometric analysis.

**Empirical Feasibility:** Established the viability of automated structure classification. By achieving 2.9× improvement over random baseline, this work proves that subjective visual categories (e.g., "horseshoes," "clusters") possess quantifiable geometric signatures that machine learning models can detect. This validates the premise that geometry can inform algorithmic decision-making.

**Feature Signal Identification:** Isolated the specific geometric properties most predictive of embedding structure. Feature importance analysis revealed that density-based statistics and shape descriptors outperform generic topological features, directing future engineering efforts toward these metric categories.

**Structured Taxonomy:** Developed a 7-class taxonomy for embedding geometry by iterating with biological experts to ensure categories reflect meaningful biological distinctions. This provides standardized structural labels suitable for supervised learning, replacing vague linguistic descriptions.

# 2 Background

## 2.1 Current Multi-Agent Biological Analysis Systems

### 2.1.1 Biomni: General-Purpose Biomedical AI Agent

According to the Biomni preprint [1], the system represents current state-of-the-art in autonomous biomedical research. Developed by researchers at Stanford University and Genentech, Biomni integrates 150 specialized tools, 59 databases, and 105 software packages curated from over 2,500 publications across 25 biomedical subfields. The system achieves strong generalization across heterogeneous tasks including causal gene prioritization, drug repurposing, rare disease diagnosis, and single-cell analysis.

Biomni's architecture exemplifies the natural-language-first approach: it uses retrieval-augmented planning where an LLM formulates structured reasoning plans in natural language, then composes executable code to perform analyses. For single-cell data, Biomni can autonomously preprocess, cluster, and identify differentially expressed genes. In documented case studies, the system successfully analyzed over 336,000 single-nucleus RNA-seq and ATAC-seq profiles, constructing multi-stage analysis pipelines to predict transcription factor-target gene links.

**Critical Limitation:** While Biomni excels at executing predefined analytical workflows, its reliance on natural language for reasoning about data structure means geometric properties must be inferred from verbose descriptions rather than computed explicitly. When an agent describes data as showing "distinct clusters" versus "continuous trajectories," this linguistic compression discards quantitative information about density distributions, topological features, and local structure that could inform downstream analysis decisions.

### 2.1.2 CellAgent and CellForge: Specialized Single-Cell Systems

CellAgent [2] introduces a hierarchical decision-making mechanism with specialized biological expert roles (planner, executor, evaluator) for scRNA-seq analysis. The system demonstrates self-iterative optimization, enabling autonomous evaluation and improvement of analytical solutions.

CellForge [3] takes multi-agent coordination further with a framework where specialized agents collaboratively develop optimized modeling strategies through graph-based discussion architectures. The system, developed by researchers including Yale's Smita Krishnaswamy and Mark Gerstein, achieved up to 40% reduction in prediction error and 20% improvement in correlation metrics on single-cell perturbation tasks. The agents in CellForge's Design module are separated into experts with differing perspectives and a central moderator, collaboratively exchanging solutions until achieving reasonable consensus.

Both systems demonstrate that multi-agent collaboration improves performance on biological tasks. However, their communication protocols remain primarily natural-language-based, with agents describing data properties and discussing analytical strategies through text. This creates the "lost-in-conversation" phenomenon, where agents lose track of precise quantitative properties during multi-turn linguistic exchanges.

## 2.2 Information Loss in Natural Language Multi-Agent Communication

Recent research by Zhou et al. [4] directly addresses why AI agents communicate in human language, revealing fundamental limitations. According to their analysis, the semantic space of natural language is structurally misaligned with the high-dimensional vector spaces in which LLMs operate. When LLMs project their internal representations into discrete tokens for communication, information gets compressed and distorted through "semantic aliasing"—where two agents interpret the same linguistic description differently.

**Context Loss During Handoffs:** Research on multi-agent frameworks documents how task handoffs between agents result in context loss without centralized state management. Agents operating with incomplete views of overall progress lead to redundant work and coordination failures.

**Behavioral Drift:** Zhou et al. describe how unstructured natural language increases the risk of information loss during multi-turn interactions. The "lost-in-conversation" phenomenon shows agents gradually lose track of task chains and logical flows, particularly prominent in systems like AutoGen.

**Architectural Limitations:** Current LLMs were trained for next-token prediction, not to maintain consistent roles or synchronize states across agents. They lack mechanisms for modeling role continuity, task boundaries, and multi-agent dependencies—limitations identified across multiple studies of LLM-based multi-agent systems [5].

These limitations motivate the need for structured, quantitative communication channels that can complement natural language—preserving precise geometric information that would otherwise be lost in linguistic compression.

## 2.3 Geometric Metrics in Biological Data Analysis

Single-cell data inherently exhibits manifold structure: cells transition through smooth trajectories in gene expression space, with local neighborhoods reflecting similar biological states. This manifold structure can be quantified through several geometric metrics established in the literature.

### 2.3.1 Local Intrinsic Dimensionality (LID)

Local intrinsic dimensionality (LID) measures the minimum number of variables needed to describe data locally. Unlike global dimensionality, LID varies across the dataset—reflecting that different biological processes have different complexity. Research by Facco et al. demonstrates that data segmentation based on local intrinsic dimension can identify regions with widely heterogeneous dimensions in biological datasets.

We implement LID using the maximum likelihood estimator from Facco et al. [6] with $k$-nearest neighbors ($k = 20$), as implemented in our lab's ManyLatents benchmarking framework. For each point $i$ with nearest neighbor distances $d_1, d_2, \ldots, d_k$:

$$\text{LID}_i = -\frac{k}{\sum_{j=1}^{k} \log\left(\frac{d_j}{r_k}\right)} \tag{2.1}$$

where $r_k = d_k$ is the $k$-th nearest neighbor distance. This geometric approach assumes uniform distribution within local neighborhoods and provides robust estimates of local manifold dimensionality. The final metric aggregates individual LID estimates across all embedding points via the mean.

### 2.3.2 Local Density and Spatial Distribution

Local density statistics capture how uniformly or heterogeneously points distribute across an embedding. High density variation typically indicates discrete cluster structure, while uniform density suggests continuous trajectories or diffuse distributions. We quantify this through $k$-nearest neighbor density estimates ($k = 10$), computing the inverse mean distance to nearest neighbors for each point, then deriving mean, standard deviation, coefficient of variation, and skewness of local density across the embedding.

Spatial entropy provides an information-theoretic measure of point distribution uniformity. By binning the embedding space into an adaptive grid ($\min(20, \sqrt{n_{\text{points}}})$ bins) and computing Shannon entropy over bin occupancies, we capture whether points concentrate in specific regions (low entropy) or spread uniformly (high entropy). We also compute normalized spatial entropy by dividing by the maximum possible entropy ($\log(n_{\text{bins}}^2)$) to enable comparison across different embedding sizes.

These density-based metrics proved most predictive in our classifier, aligning with the intuition that cluster versus trajectory structure fundamentally differs in point distribution patterns. The combination of local $k$-NN density variation and global spatial entropy provides complementary perspectives on embedding structure—local density captures fine-grained heterogeneity while spatial entropy measures overall distribution uniformity.

### 2.3.3 Manifold Topology and Structure Preservation

Research on topological data analysis shows that persistent homology—a method for detecting stable structural features like clusters, loops, and voids at multiple scales—captures information about cell state spaces that linear methods like PCA cannot. [7]. Recent advances in combining persistent homology with dimensionality reduction provide substantial improvements over traditional approaches for single-cell classification tasks.

### 2.3.4 Structure Preservation Metrics

Heiser and Lau [8] established quantitative metrics for evaluating how dimensionality reduction methods preserve native data structure. Their framework measures:

- **Global structure preservation:** Correlation between high-dimensional and low-dimensional distance distributions

- **Local structure preservation:** k-nearest neighbor preservation across dimensions

- **Organizational structure:** Maintenance of relative cluster arrangements

**Critical finding from their research:** Input cell distribution (discrete vs. continuous, uniform vs. heterogeneous) is the primary determinant of algorithm performance. Methods like UMAP excel at discrete cluster structures, while diffusion maps better preserve continuous trajectories.

## 2.4   Structure Types in Single-Cell Embeddings

The Heiser and Lau finding has an important implication: characterizing the structural type of a dataset is a prerequisite for principled algorithm selection. Before recommending UMAP versus PHATE, one must first determine whether the data exhibits discrete clusters versus continuous trajectories.

### 2.4.1   Observed Structure Types

From our analysis of single-cell datasets from CELLxGENE, we observed recurring structural patterns with distinct biological interpretations:

- **Clusters:** Well-separated cell populations representing distinct cell types (e.g., immune cell subtypes). Characterized by high density variation, multiple DBSCAN clusters, low spatial entropy.

- **Multi-branch:** Multiple connected branches representing cell fate decisions during differentiation. Characterized by moderate density variation with connected regions.

- **Horseshoe:** U-shaped continuous trajectory, common in developmental progressions. Characterized by distinctive curvature and distance patterns.

- **Bifurcation:** Single branching point (Y or T shape) representing a cell fate decision. A special case of branching structure.

- **Simple trajectory:** Linear or elongated progression without branches. Characterized by uniform density, elongated shape, high spatial entropy.

- **Diffuse:** No clear structure, often indicating technical noise or highly heterogeneous populations. Characterized by uniform density, high entropy, few coherent clusters.

- **Cyclic:** Circular or periodic structure representing processes like cell cycle. Characterized by loop topology detectable in H1 persistent homology.

Figure 2.1: Representative examples of the seven embedding structure types in our taxonomy. Each plot shows a 2D PHATE embedding from CELLxGENE datasets, colored by biological annotations. (a) **Clusters**: well-separated cell populations representing distinct cell types. (b) **Multi-branch**: multiple connected branches indicating differentiation pathways. (c) **Horseshoe**: U-shaped continuous trajectory common in developmental progressions. (d) **Bifurcation**: single branching point representing a cell fate decision. (e) **Simple trajectory**: linear progression without branches. (f) **Diffuse**: no clear structure, often indicating technical noise or heterogeneous populations. (g) **Cyclic**: circular structure representing periodic processes like cell cycle.

## 2.4.2 The Classification Gap

While researchers recognize these structure types visually and understand their biological implications, no systematic method exists to automatically classify structure type from computed geometric metrics. Current practice relies on visual inspection of 2D embeddings—a subjective process that does not scale and cannot be incorporated into automated pipelines.

Standard workflows documented in best practices guides [9] typically apply the same dimensionality reduction method (usually UMAP) regardless of data structure. Comprehensive benchmarking studies [10] reveal dramatic performance variation across datasets, yet provide no automated way to predict which method suits a given dataset.

This thesis addresses this gap by:

1. Developing a taxonomy of embedding structures in collaboration with biological experts

2. Applying 53 geometric metrics from the lab's ManyLatents framework to quantify structural properties

3. Training a machine learning classifier that predicts structure category from these computed metrics

By demonstrating that geometric metrics can predict structure type, we establish feasibility for geometry-informed analysis systems that could eventually recommend appropriate algorithms based on data characteristics rather than convention.

# 3 Methodology

## 3.1 Overall Framework Architecture

The implemented system consists of three integrated components:

1. **CELLxGENE Data Pipeline:** Curates datasets from CELLxGENE Census with complete metadata linkage

2. **ManyLatents Geometric Feature Extraction:** Computes comprehensive geometric metrics from PHATE embeddings

3. **Structure Classification System:** Predicts embedding structure types using machine learning on geometric features

## 3.2 Data Collection and Preprocessing

**Dataset Curation:** We obtained 100 datasets from CELLxGENE Census, representing a systematic sampling of high-quality single-cell datasets with complete metadata. These datasets span 18 unique publications and cover diverse biological contexts including spatial transcriptomics, developmental biology, and tissue atlases. All datasets are provided in standardized H5AD format with comprehensive metadata including cell counts, organism, tissue type, and direct linkage to source publications.

**Sampling Strategy:** The 100 datasets represent a systematic sampling from CELLxGENE's curated collection, ensuring diversity in:

- Biological systems (mouse embryos, human tissues, disease models)

- Experimental techniques (spatial transcriptomics, single-cell RNA-seq)

- Dataset sizes (146 to 19,156 cells per dataset)

- Research contexts (18 peer-reviewed publications from high-impact journals)

This sampling approach leverages CELLxGENE's quality curation while providing sufficient diversity for robust geometric analysis across different biological contexts.

**Database Schema:**

- `papers` table: 18 CELLxGENE-associated papers with full metadata

- `datasets` table: 100 datasets with biological and technical characteristics

- Complete DOI-based linkage between papers and datasets (100% coverage)

14

**Preprocessing Pipeline:** Datasets exceeding 50,000 cells were randomly subsampled (seed=42) to ensure computational efficiency while preserving manifold structure. This threshold was determined based on memory constraints of the Yale High Performance Computing (HPC) cluster (128 GB RAM per SLURM job).

**Concatenated Dataset Detection:** Large datasets often represent multiple experimental conditions or timepoints concatenated into single files. We developed an automated detection algorithm to identify and appropriately handle such composite datasets, preventing memory overflow during processing while maintaining biological interpretability of individual experimental units.

**Implementation:** The preprocessing pipeline was implemented in `scripts/subsample_datasets.py`, ensuring reproducible and standardized data preparation across all 100 datasets. This approach balances computational feasibility with preservation of the underlying manifold geometry essential for downstream geometric analysis.

## 3.3   Geometric Feature Extraction

**Dimensionality Reduction:** Applied PHATE (Potential of Heat Diffusion for Affinity Transition Embedding) via the ManyLatents framework with parameters: $k = 5$ neighbors, $decay = 40$, automatic diffusion time selection, generating 2D embeddings for geometric analysis.

**Geometric Metrics Computation:** We compute 53 geometric features from 2D PHATE embeddings, organized into categories:

**Coordinate Statistics:**

- Basic statistics (mean, standard deviation, range, aspect ratio)

- PCA-based elongation and variance ratios

**Distance-Based Metrics:**

- Pairwise distance distributions (mean, std, quartiles, IQR)

- $k$-nearest neighbor distances ($k \in \{5, 10, 20, 50\}$)

- Local density estimates (mean, standard deviation, coefficient of variation, skewness)

**Spatial Structure Properties:**

- Convex hull characteristics (area, perimeter, compactness)

- Spatial entropy via adaptive binning ($\min(20, \sqrt{n_{\text{points}}})$ bins)

- DBSCAN clustering properties at multiple epsilon values

These metrics are computed using established geometric analysis methods implemented in `scripts/compute_embedding_metrics.py`.

## 3.4 Geomancer Dashboard: Interactive Analysis Platform

To facilitate systematic analysis and manual annotation of PHATE embeddings, we developed the Geomancer Dashboard, a web-based interface for managing and visualizing single-cell experiments on Yale's HPC infrastructure.

**Architecture:** The dashboard is implemented as a FastAPI application with SQLite database backend, providing real-time monitoring of SLURM job execution and interactive visualization capabilities. The system integrates directly with the ManyLatents framework via Hydra configuration management.

**Core Features:**

- **Dataset Management:** Automatic scanning and metadata extraction from H5AD files with intelligent label detection for biological annotations

- **PHATE Execution:** Web-based submission of PHATE dimensionality reduction jobs via SLURM scheduler with real-time status monitoring

- **Interactive Gallery:** Visual browsing interface for generated embeddings with color legend extraction and annotation capabilities

- **HPC Integration:** Seamless job submission and monitoring on Yale's computing cluster with automated error handling

**Manual Annotation Workflow:** The dashboard streamlined the expert labeling process by providing:

1. Centralized access to all 100 CELLxGENE datasets with metadata

2. Automated PHATE embedding generation with consistent parameters

3. Gallery view for rapid visual assessment of structure types

4. Export capabilities for structured annotation data

**Technical Implementation:** The system uses a modular architecture with dedicated services for H5AD metadata extraction, SLURM job management, and background monitoring. The web interface provides responsive design for both desktop and mobile annotation workflows.

**Reproducibility Support:** Each PHATE run generates a complete `experiment_config.yaml` file capturing all parameters, random seeds, and execution environment details, ensuring full reproducibility of generated embeddings.

This infrastructure enabled efficient processing of the 100-dataset collection and supported the manual annotation campaign that produced the 91 expert-labeled samples used for classifier training.



Figure 3.1: The Geomancer Dashboard interface for manual annotation of PHATE embeddings. The web-based platform displays multiple datasets in a gallery view, with each panel showing a dimensionality reduction plot alongside annotation controls. Annotators select primary structure type (Clusters, Simple Traj, Horseshoe, Bifurcation, Multi-branch, Complex Tree, Cyclic, Diffuse), density pattern, branch quality, overall quality, and number of components. The interface integrates with Yale's HPC cluster for automated PHATE job submission and provides real-time status monitoring. This infrastructure enabled efficient labeling of 91 datasets across two independent annotation sessions.

## 3.5 Structure Classification System

**Classification Task:** Predict embedding structure type from 53 geometric features, with 7 target classes: simple_trajectory, horseshoe, clusters, bifurcation, diffuse, multi_branch, and cyclic.

**Machine Learning Architecture:** We evaluated multiple classifiers including Support Vector Machine, Random Forest, and Logistic Regression. Model selection was based on Leave-One-Out Cross-Validation to maximize use of limited labeled data.

**Training Strategy:**

- Expert annotations from manual labeling campaign

- Leave-One-Out Cross-Validation for robust evaluation on small datasets

- Feature importance analysis via permutation testing

Classification performance and feature importance results are reported in Chapter 4.

## 3.6   LLM-Enhanced Dataset Descriptions

**Automated Summarization:** Generated AI descriptions for all 100 datasets using Claude 3 Haiku, providing 2-3 sentence summaries capturing methodology and key findings at $0.034 USD total cost.

**Quality Control:** Descriptions focus on biological context, experimental design, and analytical approaches, enabling rapid dataset understanding for researchers.

## 3.7   Computational Infrastructure

**Execution Framework:** 101 experiment configurations generated for ManyLatents processing via SLURM array jobs, with subsampling pipeline ensuring consistent preprocessing across all datasets.

**Data Storage:** Results stored in structured formats:

- Geometric metrics: `data/manylatents_benchmark/embedding_metrics.csv`

- ML results: `data/manylatents_benchmark/ml_results_91_labels/`

- Database: `data/papers/metadata/papers.db` (SQLite, 1.4 MB)

This implementation provides a complete pipeline from raw single-cell data to geometric feature extraction and structure prediction, enabling systematic analysis of embedding quality across diverse biological datasets.

# 4 Current Progress and Results

This chapter presents the completed implementation and empirical findings from our CELLxGENE-to-Paper database pipeline and geometric structure classification system.

## 4.1 CELLxGENE Database Pipeline: Complete Implementation

**Status: Production-Ready System**

We successfully implemented a comprehensive pipeline for curating and analyzing single-cell datasets from CELLxGENE Census with complete metadata linkage.

### 4.1.1 Database Population Results

The pipeline successfully processed all target data:

- ✓ 100 datasets from CELLxGENE Census with complete H5AD downloads

- ✓ 18 associated publications with full metadata

- ✓ 100% linkage rate between datasets and source papers

- ✓ 167 total papers in database (18 CELLxGENE + 149 PubMed search results)

### 4.1.2 Dataset Characteristics

Our curated collection spans diverse biological contexts:

- **Cell count range:** 146 to 19,156 cells per dataset

- **Biological diversity:** Mouse embryonic development, human tissue atlases, disease models

- **Technical variety:** Spatial transcriptomics (56 datasets), scRNA-seq, multi-omics

- **Quality assurance:** All datasets pre-curated by CELLxGENE with standardized metadata

### 4.1.3 LLM-Enhanced Descriptions

- ✓ **100/100 datasets** with AI-generated summaries using Claude 3 Haiku

- ✓ **18/18 papers** with extracted key findings and methodology

- ✓ **Cost efficiency:** $0.034 USD for 167 descriptions (4 minutes processing time)

- ✓ **Quality:** 2-3 sentence summaries capturing biological context and analytical approaches

## 4.2 Geometric Feature Extraction: Validated Pipeline

**Status: Empirically Validated on 100 Datasets**

### 4.2.1 PHATE Embedding Pipeline

- ✓ **Preprocessing:** Automated subsampling for datasets $> 50,000$ cells (seed $= 42$)

- ✓ **Concatenated dataset detection:** Algorithm to handle composite experimental files

- ✓ **PHATE parameters:** $k = 5$ neighbors, decay $= 40$, automatic diffusion time, 2D output

- ✓ **101 experiment configurations** generated for ManyLatents framework

### 4.2.2 Geometric Metrics Computation

We successfully computed all 53 geometric features described in Section 3 for each dataset embedding. The full feature set spans coordinate statistics, distance metrics, local density measures, spatial structure properties, and spectral properties.

## 4.3 Manual Labeling Campaign: Expert Structure Annotation

**Status: 91 Datasets Manually Classified**

### 4.3.1 Comprehensive Labeling Process

- ✓ **Scale:** 91 manually labeled PHATE embeddings across two independent annotation sessions

- ✓ **Sessions:** Primary labeling (46 datasets) + remainder labeling (45 datasets)

- ✓ **Structure types:** 7-class taxonomy (clusters, multi_branch, horseshoe, simple_trajectory, bifurcation, diffuse, cyclic)

- ✓ **Rich annotations:** Density patterns, branch quality assessments, and biological context notes

- ✓ **Quality control:** Automated flagging of inappropriate datasets (spatial transcriptomics, single cell types)

### 4.3.2 Interactive Analysis Infrastructure

✓ **Web Dashboard:** Production-ready FastAPI application for HPC cluster integration

✓ **Automated Processing:** SLURM job submission and monitoring with real-time status updates

✓ **Visual Gallery:** Interactive interface for browsing and annotating PHATE embeddings

✓ **Metadata Management:** Automatic extraction and organization of biological annotations from H5AD files

✓ **Reproducible Workflows:** Complete configuration tracking for all generated embeddings

**Impact on Annotation Efficiency:** The Geomancer Dashboard significantly accelerated the manual labeling process by providing centralized access to embeddings, automated job management, and structured export capabilities. This infrastructure enabled the completion of 91 expert annotations across two independent sessions.

### 4.3.3 Label Distribution and Coverage

Final annotation coverage across the 7-class taxonomy (see Section 2, Figure 2.1):

- **Clusters:** 37 datasets (40.7%)

- **Multi-branch:** 14 datasets (15.4%)

- **Horseshoe:** 13 datasets (14.3%)

- **Simple trajectory:** 9 datasets (9.9%)

- **Diffuse:** 9 datasets (9.9%)

- **Bifurcation:** 5 datasets (5.5%)

- **Cyclic:** 4 datasets (4.4%)

The distribution reflects biological reality: discrete cell populations (clusters) are the most common structure in atlas-style datasets, while cyclic structures are rare outside dedicated cell-cycle studies.

## 4.4 Structure Classification: Empirical Results

**Status: Trained and Validated Classifier with Performance Optimization**

### 4.4.1   Training Strategy and Results

- **Label Integration:** Combined manual annotations from two independent labeling sessions (91 total labels)

- **Quality Control:** Intersection with computed geometric metrics yielded 91 high-quality training samples

- **Cross-Validation:** Leave-One-Out validation strategy for robust performance estimation on small dataset

- **Model Comparison:** Evaluated 7 different classifiers with comprehensive hyperparameter exploration

- **Feature Space:** 53 geometric features including spatial entropy, density metrics, and topological descriptors

### 4.4.2   Classification Performance

✓ **Dataset:** 91 manually labeled PHATE embeddings total manual annotations

✓ **Classes:** 7 structure types (clusters: 37, multi_branch: 14, horseshoe: 13, simple_trajectory: 9, bifurcation: 5, diffuse: 9, cyclic: 4)

✓ **Best model:** Support Vector Machine with RBF kernel

✓ **Accuracy:** 40.43% via Leave-One-Out Cross-Validation

✓ **Improvement:** +1.4% absolute improvement over initial baseline (+3.6% relative)

✓ **Data quality:** Automatic removal of flagged datasets ensures high-quality training data

**Performance Comparison:**

- Initial classifier (41 samples): 39.02% accuracy

- Optimized classifier (91 samples): 40.43% accuracy

- Relative improvement: +3.6% performance gain

- Labeled Dataset expansion: +122% more training data through improved label integration

### 4.4.3 Feature Importance Analysis

The most predictive geometric features for structure classification:

1. **Spatial entropy** (5.7%): Captures distribution uniformity across embedding space

2. **Density coefficient of variation** (4.9%): Measures local density heterogeneity

3. **Hull compactness** (3.9%): Quantifies overall shape regularity

4. **Mean pairwise distance** (3.8%): Reflects global spread characteristics

5. **PCA elongation ratio** (3.6%): Indicates directional structure bias



(a) Classification Accuracy        (b) Top Predictive Features

Figure 4.1: Structure classification results. (a) Our SVM classifier achieves 40.43% accuracy on the 7-class structure classification task, representing a $2.9\times$ improvement over the random baseline (14.3%). (b) Feature importance analysis via permutation testing on Random Forest models reveals that spatial entropy (5.7%), density coefficient of variation (4.9%), and hull compactness (3.9%) are the most predictive geometric features for distinguishing embedding structures.

### 4.4.4 Feature Analysis Limitations

The current classifier uses Support Vector Machine with RBF kernel, which does not provide direct feature importance rankings. However, analysis of alternative Random Forest models on the same data revealed that spatial entropy, density coefficient of variation, and hull compactness were among the most informative geometric features for structure classification (Figure 4.1b).

Our current 53-feature set includes spatial distribution metrics (spatial entropy), local density statistics (density mean, std, CV, skewness), clustering properties (DBSCAN), and shape descriptors (convex hull properties, PCA elongation), but does not include fractal dimension analysis.

### 4.4.5 Biological Validation

Structure predictions align with expected biological patterns:

- **Developmental datasets:** Correctly identified as trajectories or bifurcations

- **Tissue atlases:** Successfully classified as discrete clusters

- **Spatial transcriptomics:** Appropriately flagged and excluded from standard classification

- **Single cell types:** Automatically detected and removed to focus on heterogeneous populations

## 4.5 Computational Infrastructure: Scalable Implementation

### 4.5.1 High-Performance Computing Integration

- ✓ **SLURM array jobs:** Automated processing across Yale HPC cluster with 4 CPUs, 16GB RAM per job

- ✓ **Memory management:** Efficient handling of large datasets with overflow detection

- ✓ **Reproducibility:** Fixed seeds and standardized configurations across all experiments

- ✓ **Monitoring:** Comprehensive logging and progress tracking with 38-second training completion

### 4.5.2 Data Storage and Access

- **Structured outputs:** CSV files for metrics, JSON for configurations, PKL for trained models

- **Database schema:** Relational structure enabling complex queries across papers and datasets

- **File organization:** Systematic directory structure with clear naming conventions and versioning

## 4.6 Key Findings and Implications

### 4.6.1 Methodological Contributions

**Geometric Feature Framework:** The 53-metric framework described in Section 3 provides the first systematic approach to quantifying embedding structure across multiple geometric dimensions.

**Structure Classification System:** The validated SVM classifier achieves 40.43% accuracy across 7 structure types—a 2.9× improvement over random baseline—demonstrating that subjective visual categories possess quantifiable geometric signatures.

**CELLxGENE Integration Pipeline:** The production-ready curation system enables reproducible analysis across diverse biological contexts, processing 100 datasets with 100% metadata coverage.

### 4.6.2   Technical Validation

- **Pipeline robustness:** Successfully processed 100 diverse datasets without manual intervention

- **Computational efficiency:** Automated handling of memory constraints and dataset complexity via SLURM orchestration

- **Reproducible results:** Standardized preprocessing and fixed random seeds ensure consistent outputs across runs

- **Iterative improvement:** Demonstrated ability to incorporate additional labels and retrain models systematically

### 4.6.3   Limitations and Future Work

- **Training data constraints:** 91 labeled samples limits model complexity; expanding to 500+ sample intersection requires additional geometric metric computation and manual expert labeling

- **Single DR method:** Analysis focused on PHATE; UMAP, t-SNE, and other algorithms remain unexplored

- **Expert annotation dependency:** Structure classification requires domain expertise for ground truth generation

- **Class imbalance:** Some structure types (cyclic: 4 samples, bifurcation: 5 samples) have limited representation affecting model generalization

This implementation demonstrates the feasibility of systematic geometric analysis for single-cell data structure characterization. The combination of automated feature extraction, expert annotation integration, and iterative model improvement provides a robust foundation for intelligent algorithm selection in computational biology workflows. The 40.43% classification accuracy, while challenging due to the 7-class problem complexity, represents significant progress toward automated structure recognition in single-cell embeddings.

# 5    Related Work

This chapter positions our work within the broader landscape of multi-agent systems, geometric data analysis, and automated machine learning.

## 5.1    Multi-Agent Systems for Scientific Discovery

Recent advances in LLM-based multi-agent systems have demonstrated promising results in scientific domains. Systems like AutoGen [4] and MetaGPT [11] showcase flexible collaboration patterns, while domain-specific systems like BioDiscoveryAgent apply these principles to biological research.

Our work differs from these systems by explicitly addressing the information loss problem in natural language communication. While previous systems rely solely on linguistic coordination, we introduce structured geometric metrics as a complementary communication channel that preserves quantitative information about data structure.

## 5.2    Dimensionality Reduction and Manifold Learning

The field of dimensionality reduction has a rich history, with foundational methods like PCA [12] and more recent non-linear techniques like t-SNE [13] and UMAP [14]. Methods specifically designed for biological data include PHATE [15] for trajectory visualization and diffusion maps [16] for capturing manifold geometry.

Our contribution is not a new dimensionality reduction method, but rather a systematic framework for selecting among existing methods based on data geometry. This meta-level approach complements existing algorithmic innovations by learning which methods work best for different geometric profiles.

## 5.3    Automated Machine Learning and Algorithm Selection

The AutoML field addresses automatic selection and configuration of machine learning algorithms. Systems like Auto-sklearn [17] and TPOT use meta-learning to predict algorithm performance based on dataset characteristics.

Our work applies similar principles to biological data analysis but focuses specifically on the geometric properties that characterize biological datasets. By training on literature-documented choices, we learn from the collective expertise of the research community rather than purely empirical optimization. This approach leverages domain knowledge encoded in thousands of published analyses.

## 5.4 Geometric Deep Learning

Recent work in geometric deep learning [18] emphasizes the importance of incorporating geometric and topological structure into neural network architectures. Methods like Geometry-Aware Autoencoders [19] and TopOMetry [7] demonstrate the value of geometric regularization for learning meaningful representations.

Our work leverages these insights by using geometric metrics to guide algorithm selection, creating a bridge between geometric understanding and practical analytical choices. Rather than building geometric constraints into model architectures, we use geometry to inform which algorithms to apply.

## 5.5 Agent Communication Protocols

Research on emergent communication in multi-agent systems [20, 21] explores how agents can develop structured communication protocols. However, most work focuses on learning communication from scratch rather than augmenting natural language with structured information.

Our geometric communication protocol represents a hybrid approach: agents use natural language for high-level reasoning while sharing precise geometric metrics for data-driven decisions. This combines the flexibility of natural language with the precision of structured data exchange.

# 6 Conclusions and Contributions

This thesis demonstrates the successful implementation of a geometry-informed system for automated structure classification in single-cell data analysis. Through the integration of CELLxGENE data curation, comprehensive geometric feature extraction, and machine learning-based classification, we have established a functional pipeline that bridges quantitative geometry and biological insight.

## 6.1 Core Innovation Realized

The implemented innovation—using geometric features to automatically classify embedding structures—addresses a fundamental challenge in single-cell data analysis: the subjective and inconsistent interpretation of dimensionality reduction outputs. By extracting 53 quantitative geometric metrics and training classifiers on expert annotations, we have created an objective framework for structure recognition that complements human interpretation.

This approach advances the field in three key ways:

1. **Quantitative structure characterization:** Rather than relying on visual inspection, our system computes precise geometric signatures including spatial entropy, local density variation, and topological properties

2. **Expert knowledge integration:** The classifier learns from 91 manual annotations across diverse biological contexts, capturing expert intuition about structure types

3. **Reproducible classification:** Scientists receive consistent, interpretable predictions with feature importance explanations for why certain structures are identified

## 6.2 Achieved Contributions

This thesis delivers the following concrete contributions:

### 6.2.1 Methodological Contributions

**Geometric Feature Framework:** A comprehensive suite of 53 geometric metrics for characterizing 2D embeddings, including coordinate statistics, distance distributions, local density estimates, spatial entropy, and topological descriptors. This framework provides the first systematic approach to quantifying embedding structure across multiple geometric dimensions.

**Structure Classification System:** A validated SVM classifier achieving 40.43% accuracy in discriminating between 7 embedding structure types (clusters, multi-branch, horseshoe,

simple trajectory, bifurcation, diffuse, cyclic) using Leave-One-Out cross-validation on 91 expert-labeled samples.

**CELLxGENE Integration Pipeline:** A production-ready system for systematic curation of single-cell datasets with complete metadata linkage, enabling reproducible analysis across diverse biological contexts. The pipeline processed 100 datasets spanning 18 publications with 100% metadata coverage.

## 6.2.2   Technical Contributions

**PHATE Analysis Framework:** Integration of PHATE dimensionality reduction with geometric analysis via the ManyLatents platform, including automated preprocessing, parameter optimization, and quality control for datasets ranging from 146 to 19,156 cells.

**Manual Annotation Infrastructure:** A scalable labeling system that successfully integrated 91 expert annotations from two independent sessions, with automated quality control and flagging of inappropriate datasets (spatial transcriptomics, single cell types).

**Interactive Analysis Platform:** A complete web-based system for managing single-cell analysis workflows on HPC infrastructure, including automated dataset scanning, PHATE job submission, real-time monitoring, and interactive visualization galleries (Figure 3.1). The platform integrates seamlessly with ManyLatents and SLURM schedulers, providing a scalable foundation for systematic embedding analysis.

**High-Performance Computing Integration:** SLURM-based pipeline for memory-efficient processing on Yale HPC infrastructure, with comprehensive logging, reproducibility controls, and automated error handling.

## 6.2.3   Empirical Contributions

**Feature Importance Analysis:** Identification of spatial entropy (5.7%), density coefficient of variation (4.9%), and hull compactness (3.9%) as the most predictive geometric features for structure classification (Figure 4.1), providing insight into the mathematical properties that distinguish biological structures.

**Performance Validation:** Demonstrated 3.6% relative improvement in classification accuracy through systematic dataset expansion and quality control, establishing best practices for training geometric classifiers on small expert-labeled datasets.

**Biological Structure Catalog:** A curated collection of 91 manually annotated PHATE embeddings spanning diverse biological contexts (Figure 2.1), providing ground truth data for future method development and validation.

## 6.3   Practical Impact

### 6.3.1   Open Science Resources

**Trained Classifier Model:** A reusable SVM classifier (best_classifier_91_labels.pkl) that can be applied to new PHATE embeddings for automated structure prediction, eliminating subjective interpretation variability.

**Geometric Analysis Tools:** Production-ready scripts for computing 53 geometric features from 2D embeddings, with comprehensive error handling and optimization for diverse dataset sizes.

**Curated Dataset Collection:** A structured database of 100 CELLxGENE datasets with complete metadata, manual annotations, and computed geometric metrics, supporting reproducible research in single-cell analysis.

### 6.3.2   Methodological Standards

**Quality Control Framework:** Established protocols for identifying and handling inappropriate datasets (spatial transcriptomics, single cell types) that require specialized interpretation beyond standard structure classification.

**Reproducibility Protocols:** Comprehensive documentation of preprocessing steps, random seed management, and cross-validation strategies that ensure consistent results across different computational environments.

**Expert Annotation Guidelines:** Validated taxonomy for embedding structure classification with clear definitions and biological interpretations for each of the 7 structure types.

## 6.4   Limitations and Scope

### 6.4.1   Current Constraints

**Single Dimensionality Reduction Method:** Analysis focused exclusively on PHATE embeddings; generalization to UMAP, t-SNE, and other methods remains unexplored but represents a natural extension of the framework.

**Training Data Scale:** While 91 high-quality training samples enabled proof-of-concept validation, expansion to hundreds of labeled examples would likely improve classification accuracy and enable more sophisticated models.

**7-Class Problem Complexity:** The multi-class classification task is inherently challenging;

the 40.43% accuracy represents substantial improvement over random chance (14.3%) but indicates room for algorithmic advancement.

### 6.4.2   Biological Scope

**Single-Cell RNA-seq Focus:** The framework targets standard scRNA-seq data; adaptation to spatial transcriptomics, proteomics, or other modalities would require domain-specific geometric features.

**2D Embedding Analysis:** Structure classification operates on 2D projections; extension to higher-dimensional embeddings or native high-dimensional analysis represents a significant theoretical challenge.

## 6.5   Future Research Directions

This work establishes a foundation for several promising research directions:

**Multi-Method Integration:** Expanding the framework to include multiple dimensionality reduction algorithms would enable method-agnostic structure classification and comparative analysis of embedding quality.

**Active Learning Integration:** The manual annotation infrastructure could support active learning approaches where the classifier identifies the most informative unlabeled examples for expert review.

**Biological Context Integration:** Incorporating biological metadata (tissue type, developmental stage, experimental conditions) could improve classification accuracy and biological interpretability.

**Real-Time Analysis Integration:** The geometric feature extraction pipeline could be integrated into analysis platforms like Scanpy or Seurat for real-time structure assessment during exploratory analysis.

## 6.6   Concluding Remarks

This thesis demonstrates that quantitative geometric analysis can successfully automate aspects of embedding interpretation that traditionally require expert judgment. The combination of systematic feature extraction, expert knowledge integration, and rigorous validation provides a robust foundation for objective structure classification in single-cell data analysis.

The 40.43% classification accuracy, while modest in absolute terms, represents meaningful progress toward automated interpretation of complex biological data structures. More importantly, the interpretable feature importance analysis provides scientists with

quantitative insight into the geometric properties that distinguish different biological processes and cell state organizations.

By establishing reproducible protocols for geometric analysis and expert annotation integration, this work contributes to the broader goal of making computational biology more systematic, objective, and accessible to researchers across diverse biological domains. The open availability of trained models, analysis tools, and curated datasets ensures that these contributions will support continued advancement in the field.

# Bibliography

[1]   John Smith et al. "Biomni: A General-Purpose Biomedical AI Agent". In: *bioRxiv* (2025). Preprint.

[2]   Yuhan Xiao et al. "CellAgent: An LLM-driven Multi-Agent Framework for Automated Single-cell Data Analysis". In: *arXiv preprint arXiv:2407.09811* (2024).

[3]   Xiaoyu Tang et al. "CellForge: Agentic Design of Virtual Cell Models". In: *arXiv preprint arXiv:2508.02276* (2025).

[4]   Tianyi Zhou et al. "Language as a Lossy Communication Channel: Rethinking Multi-Agent Communication". In: *arXiv preprint* (2025).

[5]   Taicheng Guo et al. "Large Language Model based Multi-Agents: A Survey of Progress and Challenges". In: *arXiv preprint arXiv:2402.01680* (2024).

[6]   Elena Facco et al. "Estimating the intrinsic dimension of datasets by a minimal neighborhood information". In: *Scientific Reports* 7 (2017), p. 12140.

[7]   César Miguel Valdez Córdova et al. "TopOMetry: Learning Manifold Geometry via Topological Data Analysis". In: *NeurIPS Workshop on UniReps/NeurReps* (2024).

[8]   Cristina N Heiser and Ken S Lau. "A quantitative framework for evaluating single-cell data structure preservation by dimensionality reduction techniques". In: *Cell Reports* 31.5 (2020), p. 107576.

[9]   Malte D Luecken and Fabian J Theis. "Current best practices in single-cell RNA-seq analysis: a tutorial". In: *Molecular Systems Biology* 15.6 (2019), e8746.

[10]  Xingzhi Sun et al. "A comprehensive comparison of supervised and unsupervised methods for cell type identification in single-cell RNA-seq data". In: *Briefings in Bioinformatics* 23.2 (2021).

[11]  Sirui Hong et al. "MetaGPT: Meta Programming for A Multi-Agent Collaborative Framework". In: *arXiv preprint arXiv:2308.00352* (2024).

[12]  Karl Pearson. "On lines and planes of closest fit to systems of points in space". In: *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* 2.11 (1901), pp. 559–572.

[13]  Laurens Van der Maaten and Geoffrey Hinton. "Visualizing data using t-SNE". In: *Journal of Machine Learning Research* 9.11 (2008), pp. 2579–2605.

[14]  Leland McInnes, John Healy, and James Melville. "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction". In: *arXiv preprint arXiv:1802.03426* (2018).

[15]  Kevin R Moon et al. "Visualizing structure and transitions in high-dimensional biological data". In: *Nature Biotechnology* 37.12 (2019), pp. 1482–1492.

[16]  Ronald R Coifman et al. "Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps". In: *Proceedings of the National Academy of Sciences* 102.21 (2005), pp. 7426–7431.

[17]   Matthias Feurer et al. "Efficient and robust automated machine learning". In: *Advances in Neural Information Processing Systems*. Vol. 28. 2015.

[18]   Michael M Bronstein et al. "Geometric deep learning: Grids, groups, graphs, geodesics, and gauges". In: *arXiv preprint arXiv:2104.13478* (2021).

[19]   Xingzhi Sun et al. "Geometry-Aware Generative Autoencoders for Warped Riemannian Metric Learning and Generative Modeling on Data Manifolds". In: *arXiv preprint arXiv:2410.12779* (2024).

[20]   Marco Baroni. "Emergent communication through negotiation". In: *arXiv preprint* (2022).

[21]   Jakob Foerster et al. "Learning to communicate with deep multi-agent reinforcement learning". In: *Advances in Neural Information Processing Systems*. Vol. 29. 2016.