# scSSL-Bench: Benchmarking Self-Supervised Learning for Single-Cell Data

**Olga Ovcharenko** TU Berlin          Florian Barkmann ETH          Philip Toma ETH
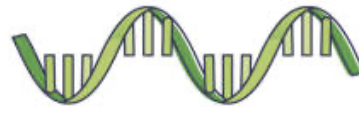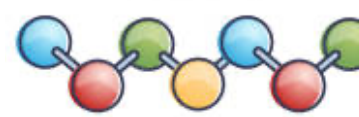
Imant Daunhawer ETH          Julia Vogt ETH          Sebastian Schelter TU Berlin          Valentina Boeva ETH
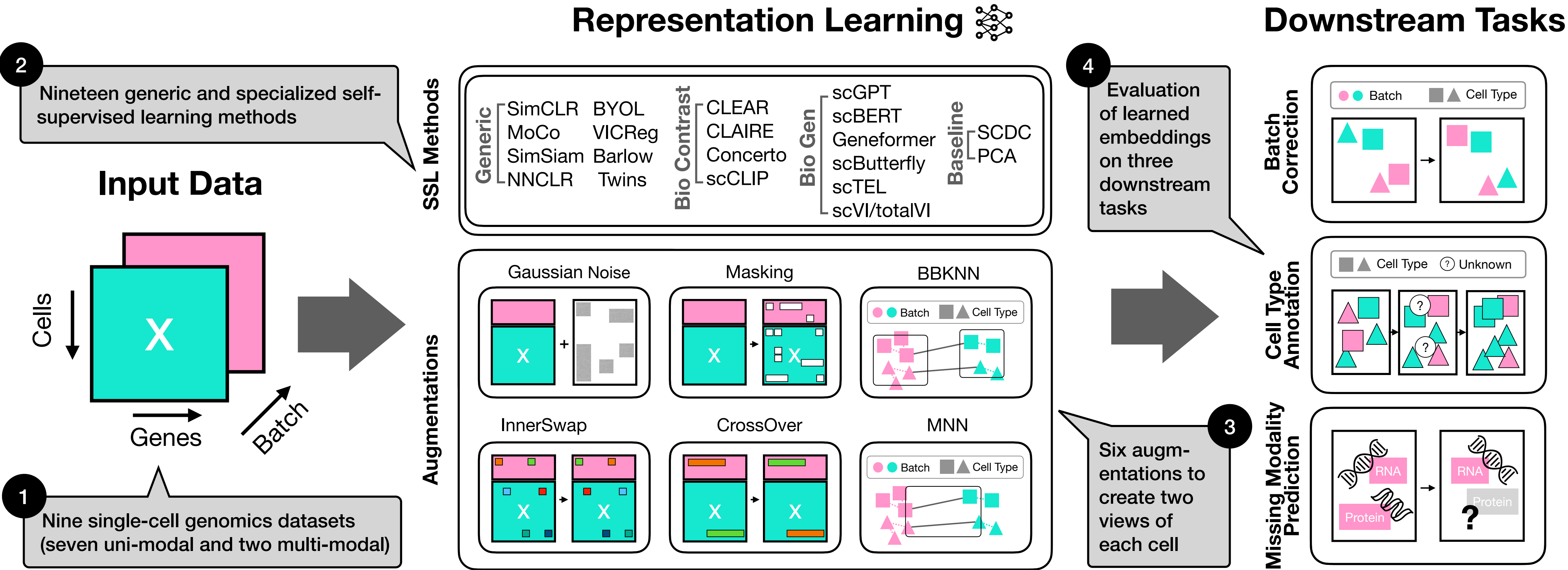
BIFOLD          TU berlin          ETH zürich

# Motivation

- **Why** is genomics analysis an important research direction?

  - Profile cells at different resolutions and modalities

  - Understand diseases, develop personalized treatments, trace origins of conditions like cancer and autoimmune disorders

- **What** is single-cell data?

  - High-dimensional gene expression (**GEX**) levels in individual cells

  - Additional molecular features measured together with **GEX**

    - **Protein** levels +

    - Open **chromatin** accessibility +

# Why is a Benchmark Needed?

- Self-supervised learning (SSL) is a powerful approach for the representation extraction from single-cell data

- **Research questions**:

  - Benchmark if **specialized** single-cell SSL methods outperform **generic** methods

  - Assess **hyperparameters** and **augmentation** techniques of generic SSL approaches

  - Evaluate if genomics benefit from techniques proposed for images

# scSSL-Bench Design



ICML 2025. Olga Ovcharenko

3

# Downstream Tasks on Learned Cell Representations

- **Batch effect correction**
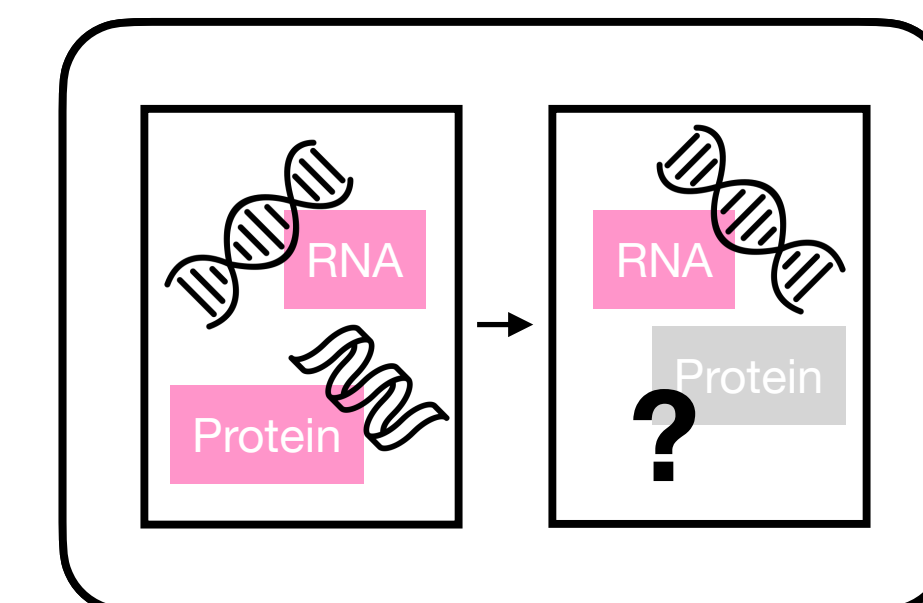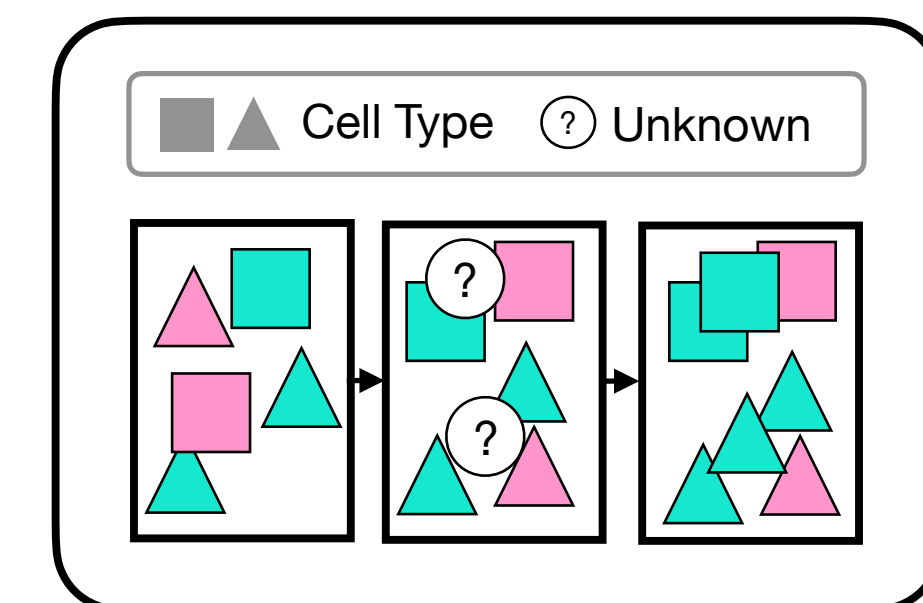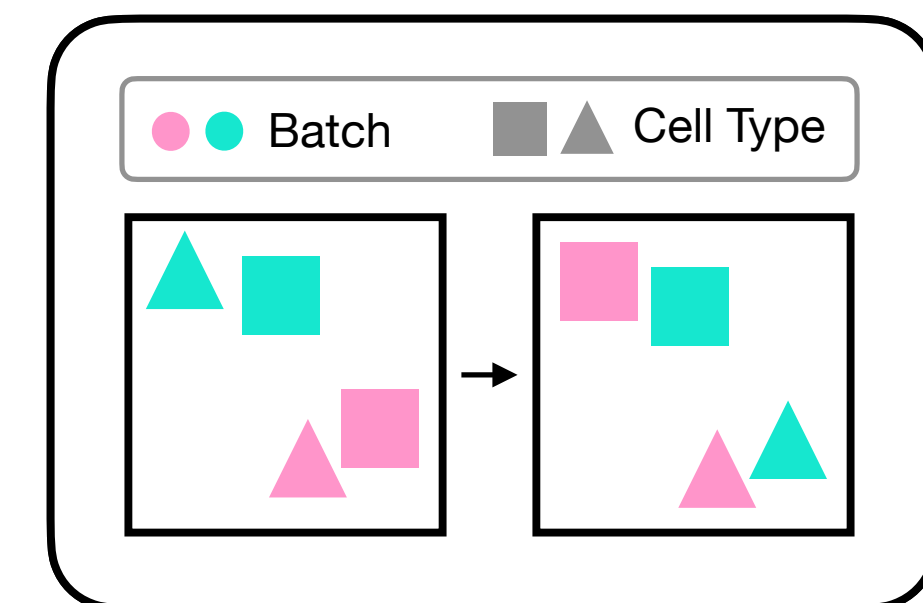
  - Reduce technical biases introduced while sequencing

  - Preserve true biological signal

- **Cell typing**

  - Annot... ...set by mapping them t...

- **Missing**

  - Infer u... ...the hold-out dataset...



Evaluation of learned embeddings on three downstream tasks

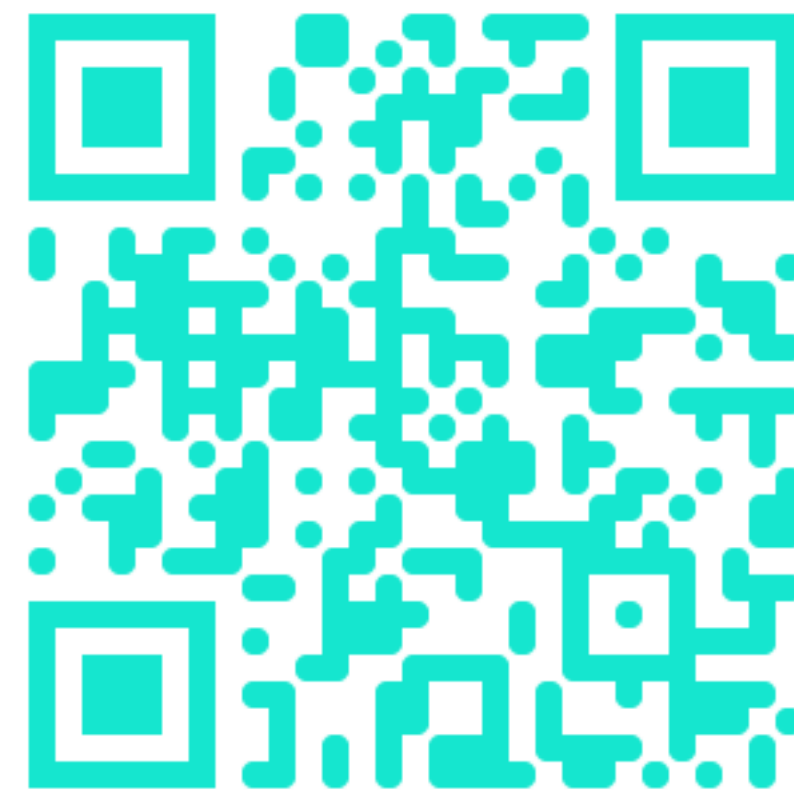Six augm- entations to create two views of each cell

# Conclusions

- **Specialized** single-cell SSL methods work better for **uni-modal** data

  - scVI, CLAIRE, fine-tuned scGPT

- **Generic** SSL methods succeed in **multi-modal** single-cell data integration

  - SimCLR, VICReg

- **Masking** is the best augmentation technique

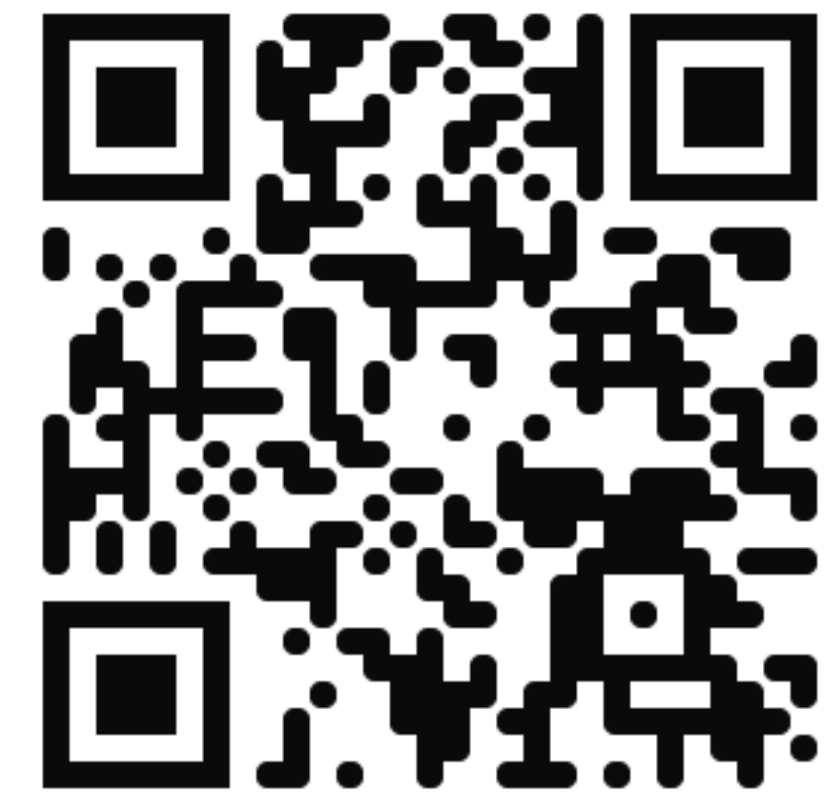- **Moderately-sized embeddings** lead to better results

# scSSL-Bench: Benchmarking Self-Supervised Learning for Single-Cell Data



**Code**

**Paper**

**Contact Me**