#### **Uncertainty Estimation for Single-cell Label Transfer Robin Khatri and Stefan Bonn**

COPA, August 25 2022



Center for Molecular Neurobiology Hamburg 



Universitätsklinikum Hamburg-Eppendorf



#### **Single-cell Transcriptomics** A brief introduction

single-cell level

#### Single-cell RNA sequencing (scRNA-seq) measures mRNA expression at a

### **Single-cell Transcriptomics A brief introduction**

single-cell level



Source: 10X Genomics

#### Single-cell RNA sequencing (scRNA-seq) measures mRNA expression at a

### **Single-cell Transcriptomics A brief introduction**

- Cells can be grouped into groups known as cell-types that share similar features and lineages
- > 2300 cell-types (Osumi-Sutherland et. al., 2021, Nature Cell Biology)



#### Cell type identification Workflow

#### Classical approach

Clustering (PCA, UMAP, tSNE)



### **Cell type identification** Workflow

#### <u>Classical approach</u>

Clustering (PCA, UMAP, tSNE)



- Dependence on clustering algorithm, Requires very good (and few) marker genes per cell type, not suitable for identifying rare cell types.

+ Detailed control, expert-based

### **Cell type identification** Label transfer

Idea: Use already available and annotated datasets

map.org) etc.

- + Machine Learning algorithms are getting better
- No (human) expert

quality issues, partial references, dynamic nature of cell states

- + Larger and larger reference datasets are being generated e.g. Human Cell Atlas (https://www.humancellatlas.org), Allen Brain Atlas (https://portal.brain-

- Many sources of error - Batch effects, data integration (batch correction), data

#### **Cell type identification** Label transfer



Illustration of Batch Effect on two peripheral blood mononuclear cell (PBMC) datasets.



Post-batch correction using Harmony (Karunsky et. al., 2019, Nature Methods)

Why use conformal prediction?

- Model agnostic
- Interpretable
- Simple and general

Goals: Explore conformal prediction, investigate how it could help in single-cell label transfer and how it could identify previously unseen cell types





Inductive conformal classifiers (ICPs) **Base classifier: SVM** 





Inductive conformal classifiers (ICPs) **Base classifier: SVM** 

Comparison: SVM, and SVM (thr = 0.7) Conformal-credibility setup

Harmony embeddings

**Conformal procedure** 

50 principal components







PBMCs	Dataset I PBMC 6k 7 PBMC 8k 9	B CD 704 2,24 992 1,97	<b>4T Cl</b> 0 714 5 89	<b>D8T Moi</b> 4 1,39' 0 1,87	nocytes 7 )	NK 7 2,240 5 320 6	<b>Fotal</b> 5,356 5,047			
	Dataset	Alpha	Beta	Gamma	Delta	Acinar	Ductal	Endo.	Stellate	Total
Pancreas	inDrop	$2,\!249$	$3,\!048$	260	613	272	898	689	362	$8,\!391$
	SS2	$1,\!109$	796	219	142	103	462	67	63	$2,\!961$
	CEL-Seq 2	885	600	125	199	170	315	31	101	$2,\!426$
	Fluidigm C1	241	300	12	21	6	27	12	13	632
	CEL-Seq	220	341	21	66	162	402	37	22	$1,\!271$

Inductive conformal classifiers (ICPs) **Base classifier: SVM** 

Comparison: SVM, and SVM (thr = 0.7) Conformal-credibility setup

Harmony embeddings

**Conformal procedure** 

50 principal components







Before batch correction

SVM performance before and after batch correction with Harmony Training set: PBMC 6k, Test set: PBMC 8k



After batch correction

С

$\mathbf{Method}$	
SVM	

SVM SVM (thr 0.7) SVM (thr 0.7) ICP (With SVM) ICP (With SVM)

#### Test set PBMC 8k PBMC 6k PBMC 8k PBMC 6k PBMC 8k

PBMC 6k

Peformance on PBMCs ICP at significance: 0.025

Average accuracy	<b>Overall accuracy</b>
0.771	0.844
0.616	0.834
0.842	0.868
0.742	0.863
0.935	0.854
0.790	0.868





Rate of prediction sets on PBMC 8k



Rate of prediction sets on PBMC 8k



Average credibility of each cell





Average credibility of each cell



le	$\mathbf{Test} \ \mathbf{set}$	cell type	Average credibility
Taalla	PBMC 8k	Bcells	0.417
TCells	PBMC 8k	CD4Tcells	0.415
Icells	PBMC 8k	CD8Tcells	0.143
ocytes	PBMC 8k	NK	0.280
	PBMC 8k	Monocytes	0.419

Average credibility of each cell type

PBM	С	8k
-----	---	----

Method	Cell type
SVM (thr 0.7)	Bcells
SVM (thr $0.7$ )	CD4Tcells
SVM (thr 0.7)	CD8Tcells
SVM (thr 0.7)	NK
SVM (thr 0.7)	Monocytes
CC-ICP	Bcells
CC-ICP	CD4Tcells
CC-ICP	CD8Tcells
CC-ICP	NK
CC-ICP	Monocytes

Desired error rate = 0.025 % per cell type

<b>Classification</b> rate	Accuracy
0.982	1
0.956	0.988
0.808	0.068
0.94	0.987
0.968	0.975
0.968	0.99
0.564	0.98
0.11	0.946
0.638	0.931
0.649	0.99



Pancreas datasets after batch correction



Pancreas datasets after batch correction



Desired error rates vs classification and accuracy

Unknow cell type	Assignment on Unknown	Average assignments on known
Alpha	0.08	0.42
Delta	0.59	0.64
Beta	0.68	0.82
Ductal	0.514	0.762
Acinar	0.501	0.761
Gamma	0.455	0.636
Endothelial	0.445	0.748
Stellate	0.379	0.705

Assignments of held-out (unknown) cell types

#### Results Impact of batch correction

Batch transfer may itself cause error rather than the model itself Evaluated two other algorithm on PBMCs:

- Scanorama (Hie et. al., 2019, *Nature Biotechnology*) lacksquare
- Scgen (Lotfollahi et. al., 2019, Nature Methods) ullet

#### Results Impact of batch correction

Batch transfer may itself cause error rather than the model itself Evaluated two other algorithm on PBMCs:

- Scanorama (Hie et. al., 2019, *Nature Biotechnology*)  $\bullet$
- Scgen (Lotfollahi et. al., 2019, Nature Methods)

Ho
0.61
0.81
0.79

Homogeneity scores for integrating PBMC datasets

#### mogeneity score

13

12

95

#### **Results** Impact of batch correction



Test set error rates Batch correction: Harmony

Should we use Scanorama always? Perhaps not. Quality of batch correction methods is data-dependent.



Test set error rates Batch correction: Scanorama

- Conformal-credibility setup provides an approach for accurate label transfer
- Uncertainties drastically affected by the batch correction method in use

- Conformal-credibility setup provides an approach for accurate label transfer
- Uncertainties drastically affected by the batch correction method in use
- Progress to be made on three ends:
  - Technical:
    - Evaluating of a set of NCMs and further development of algorithms
    - Extensive benchmarking with multiple batch correction methods (Ongoing)
    - Conformal anomaly detection for discovering unknown cell types

- Conformal-credibility setup provides an approach for accurate label transfer
- Uncertainties drastically affected by the batch correction method in use
- Progress to be made on three ends:
  - Technical:
    - Evaluating of a set of NCMs and further development of algorithms
    - Extensive benchmarking with multiple batch correction methods (Ongoing)
    - Conformal anomaly detection for discovering unknown cell types
  - Biological:
    - Going deeper into cellular subsets (*Ongoing*)
    - Discovering cell states

- Conformal-credibility setup provides an approach for accurate label transfer
- Uncertainties drastically affected by the batch correction method in use
- Progress to be made on three ends:
  - Technical:
    - Evaluating of a set of NCMs and further development of algorithms
    - Extensive benchmarking with multiple batch correction methods (Ongoing)
    - Conformal anomaly detection for discovering unknown cell types
  - Biological:
    - Going deeper into cellular subsets (Ongoing)
    - Discovering cell states
  - Datasets:

• Use atlas-level datasets (Some exciting recent developments and ongoing global efforts)

# **Recent developments**



C. DOMÍNGUEZ CONDE et. al., 2022, Science

Classification without the use of batch correction (https://www.celltypist.org), thanks to curation of vast amounts of cell types and careful harmonization of multiple public datasets