

# EFFICIENT GROUPING METHODS FOR THE ANNOTATION AND SORTING OF SINGLE CELLS

## DISPUTATION

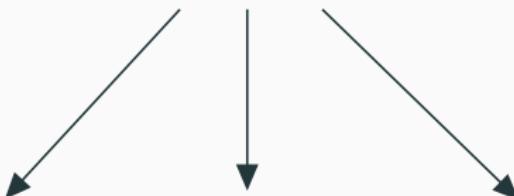
---

Markus Lux

April 24, 2018

# OVERVIEW

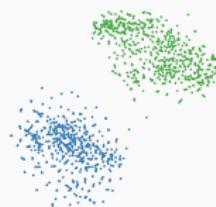
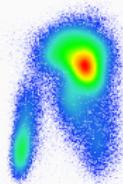
Efficient grouping methods for  
the annotation and sorting of single cells



Flow cytometry

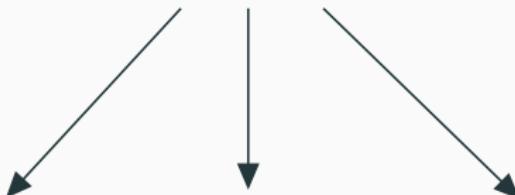
Metagenomics

Single-cell genomics

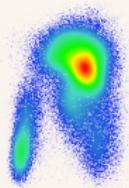


# OVERVIEW

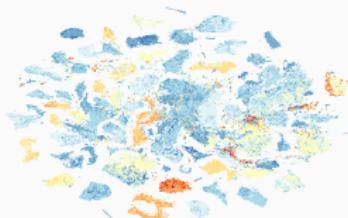
Efficient grouping methods for  
the annotation and sorting of single cells



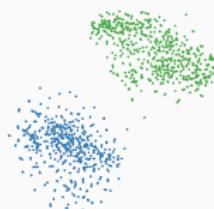
Flow cytometry



Metagenomics



Single-cell genomics



# FLOW CYTOMETRY IN A NUTSHELL

---

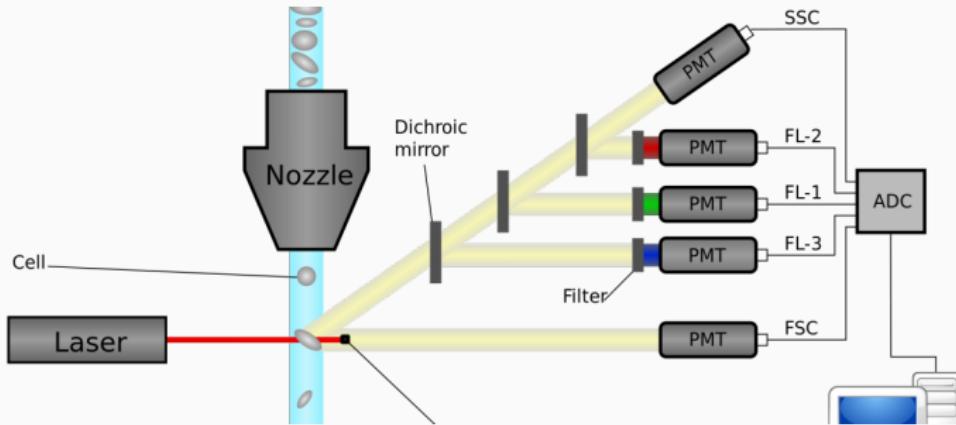
- Use laser to measure cell parameters
- Cell counting and sorting using fluorescent markers
  - Viability
  - Antigens
  - Protein expression
  - DNA characteristics

# FLOW CYTOMETRY IN A NUTSHELL

---

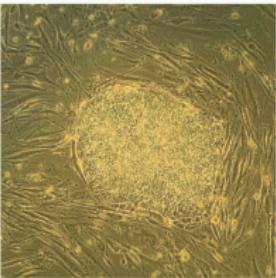
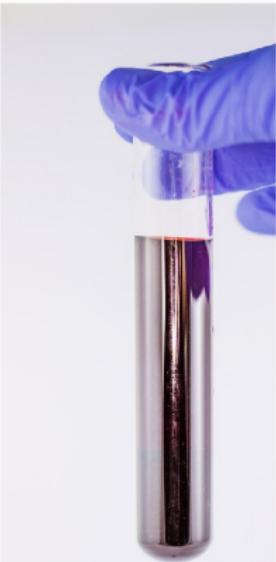
- Use laser to measure cell parameters
- Cell counting and sorting using fluorescent markers
  - Viability
  - Antigens
  - Protein expression
  - DNA characteristics

# FLOW CYTOMETRY DATA GENERATION

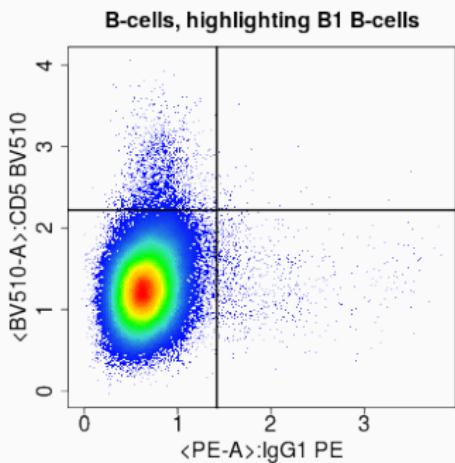


→ cells represented using  
max. 20 features (channels/markers)

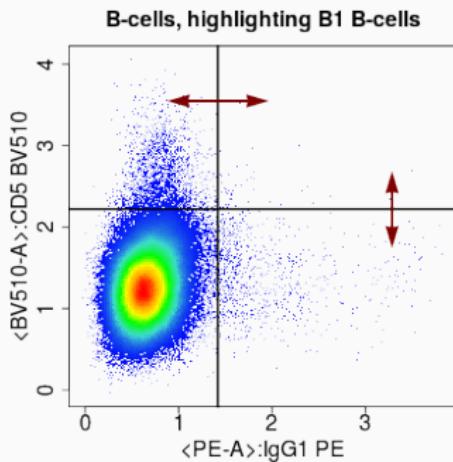
# FLOW CYTOMETRY APPLICATIONS



# GATING



# GATING



Measure of interest: **cell frequency**

# EXISTING GATING SOLUTIONS

---

- **Manual:** resource-intensive, subjective

Finak et al. (2016)

- **Clustering:** inaccurate for complex, rare populations

Aghaeepour et al. (2011); Qiu et al. (2011); Van Gassen et al. (2015); Ge and Sealfon (2012); ...

- **Supervised:** difficult to tune parameters

Finak et al. (2014); Malek et al. (2015); Li et al. (2017); ...

- **No gating:** missing acceptance, limited interpretability

Mair et al. (2016)

# EXISTING GATING SOLUTIONS

---

- **Manual:** resource-intensive, subjective  
Finak et al. (2016)
- **Clustering:** inaccurate for complex, rare populations  
Aghaeepour et al. (2011); Qiu et al. (2011); Van Gassen et al. (2015); Ge and Sealfon (2012); ...
- **Supervised:** difficult to tune parameters  
Finak et al. (2014); Malek et al. (2015); Li et al. (2017); ...
- **No gating:** missing acceptance, limited interpretability  
Mair et al. (2016)

## EXISTING GATING SOLUTIONS

- **Manual:** resource-intensive, subjective  
Finak et al. (2016)
- **Clustering:** inaccurate for complex, rare populations  
Aghaeepour et al. (2011); Qiu et al. (2011); Van Gassen et al. (2015); Ge and Sealfon (2012); ...
- **Supervised:** difficult to tune parameters  
Finak et al. (2014); Malek et al. (2015); Li et al. (2017); ...
- **No gating:** missing acceptance, limited interpretability  
Mair et al. (2016)

## EXISTING GATING SOLUTIONS

- **Manual:** resource-intensive, subjective  
Finak et al. (2016)
- **Clustering:** inaccurate for complex, rare populations  
Aghaeepour et al. (2011); Qiu et al. (2011); Van Gassen et al. (2015); Ge and Sealfon (2012); ...
- **Supervised:** difficult to tune parameters  
Finak et al. (2014); Malek et al. (2015); Li et al. (2017); ...
- **No gating:** missing acceptance, limited interpretability  
Mair et al. (2016)

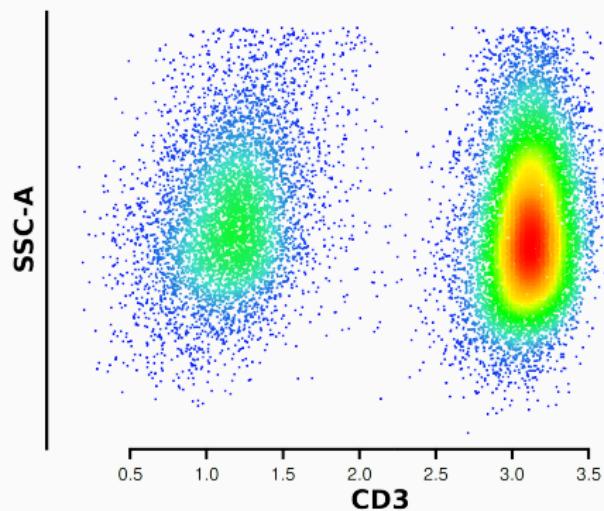
## EXISTING GATING SOLUTIONS

- **Manual:** resource-intensive, subjective  
Finak et al. (2016)
- **Clustering:** inaccurate for complex, rare populations  
Aghaeepour et al. (2011); Qiu et al. (2011); Van Gassen et al. (2015); Ge and Sealfon (2012); ...
- **Supervised:** difficult to tune parameters  
Finak et al. (2014); Malek et al. (2015); Li et al. (2017); ...
- **No gating:** missing acceptance, limited interpretability  
Mair et al. (2016)

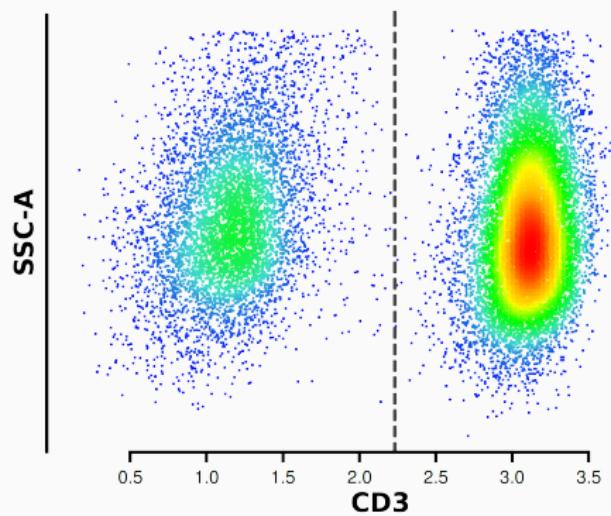
→ Semi-supervised approach:  
**flowLearn**

Lux et al. (2018)

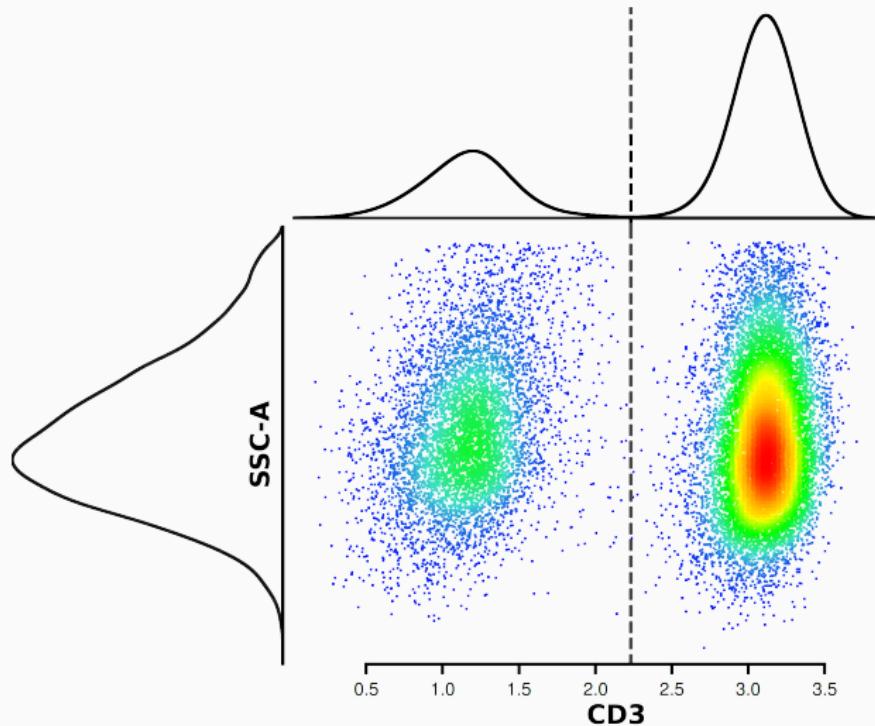
# FLOWLEARN APPROACH



## FLOWLEARN APPROACH

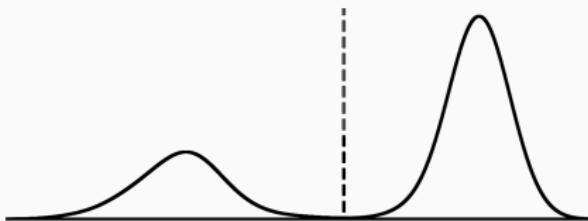


## FLOWLEARN APPROACH

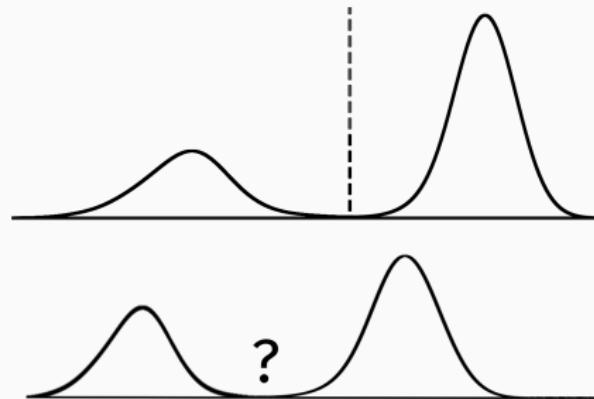


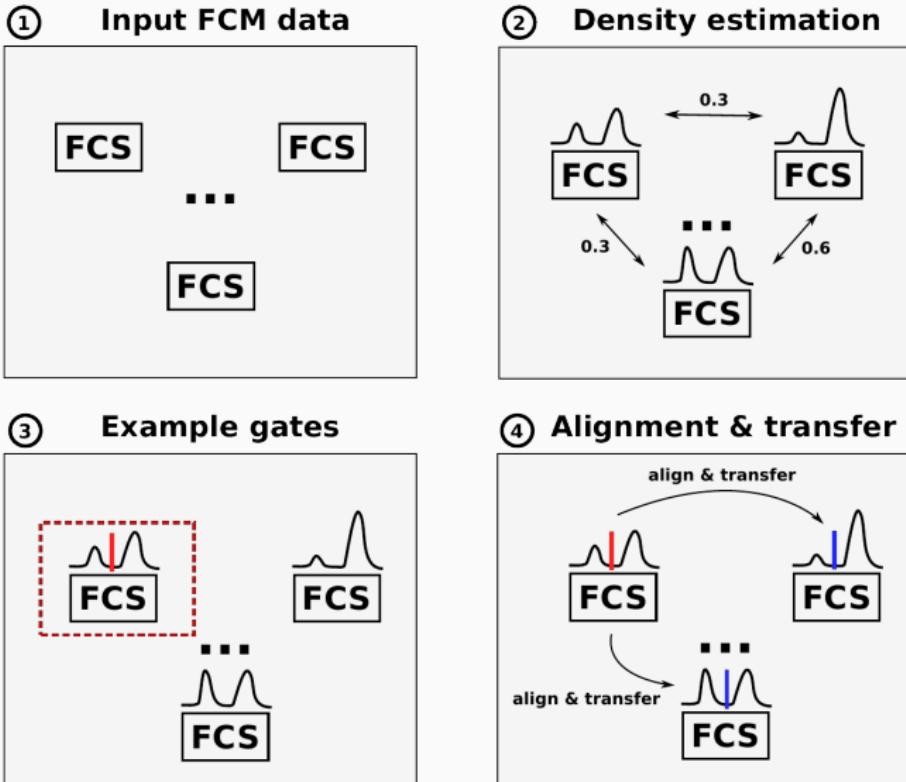
# FLOWLEARN APPROACH

---

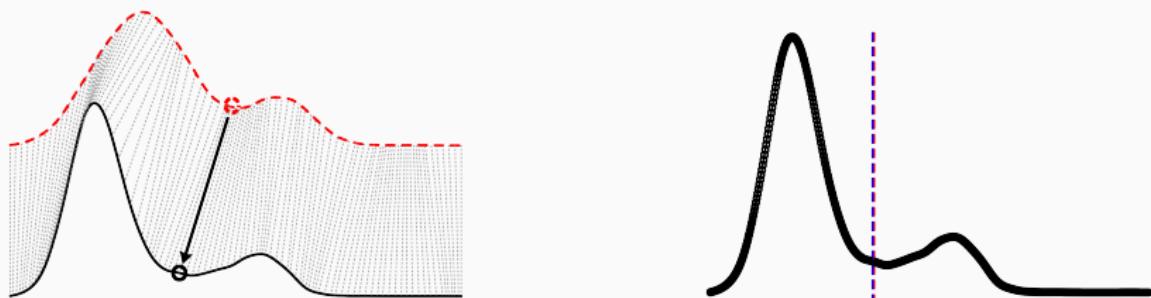


# FLOWLEARN APPROACH



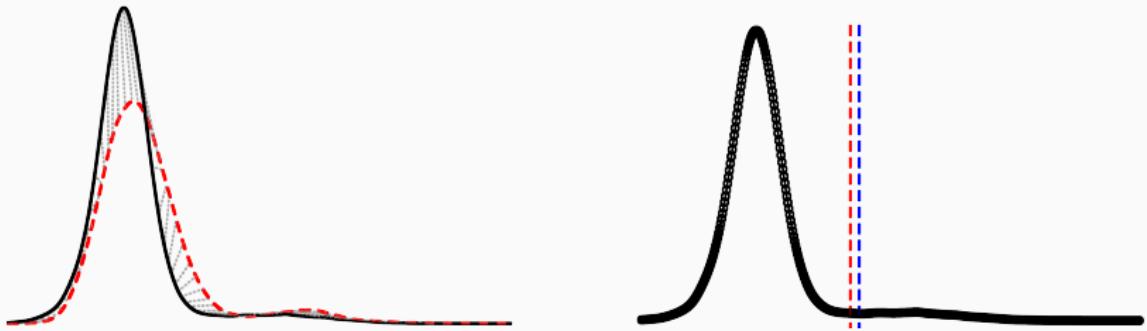


# FLOWLEARN: ALIGNMENT AND PREDICTION



1. Kernel density estimation & Smoothing splines
2. Derivative Dynamic Time Warping

## GATING RARE POPULATIONS



Rare populations extend density tail.

## Benchmark data sets

- Mice (1 dataset, 2665 samples)
- FlowCAP (4 datasets, 63 samples, 7 centers)

## Benchmark data sets

- Mice (1 dataset, 2665 samples)
- FlowCAP (4 datasets, 63 samples, 7 centers)

## Metrics

- per population
- $F_1$ -score
- cell frequency error

## RESULTS

---

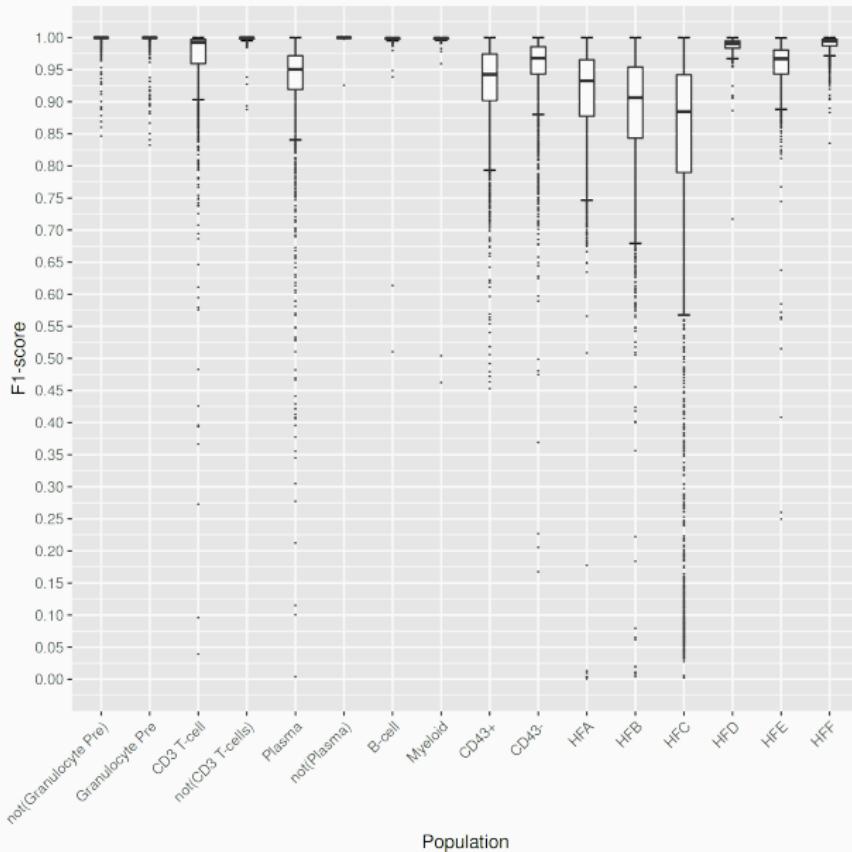
	# populations	$F_1 > 0.9$	$F_1 > 0.99$
Mice	16	15	9
FlowCap	48	36	11
Total	64	51	20

## RESULTS

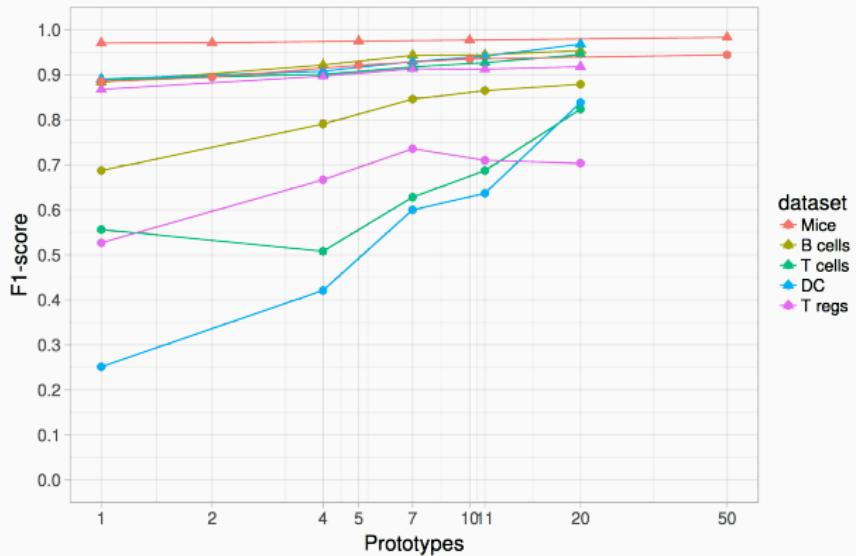
	# populations	$F_1 > 0.9$	$F_1 > 0.99$
Mice	16	15	9
FlowCap	48	36	11
Total	64	51	20

- Significantly better than current state-of-the-art methods
  - FlowSOM, Van Gassen et al. (2015)
  - DeepCyTOF, Li et al. (2017)

# MICE, 1 PROTOTYPE

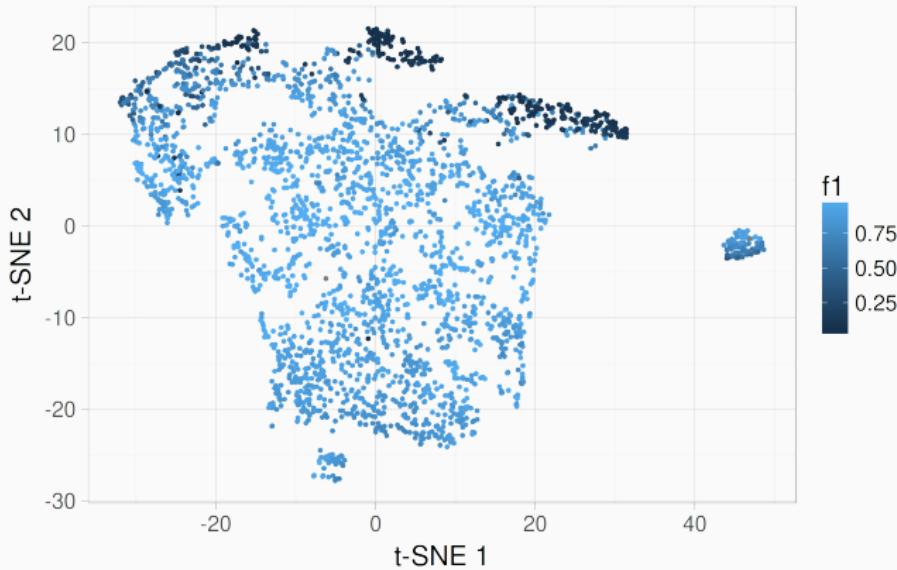


# INFLUENCE OF THE NUMBER OF PROTOTYPES



## QUALITY CHECKING

Use flowLearn for checking quality of existing gates.



## OUTLOOK

---

- Investigate other prototype choices
- Extend alignment to two dimensions
- Collaboration with clinical research

MARKUS LUX, RYAN REMY BRINKMAN, CEDRIC CHAUVE, ADAM LAING, ANNA LORENC, LUCIE ABELER-DÖRNER, BARBARA HAMMER

*FLOWLEARN: FAST AND PRECISE IDENTIFICATION AND QUALITY CHECKING OF CELL POPULATIONS IN FLOW CYTOMETRY*

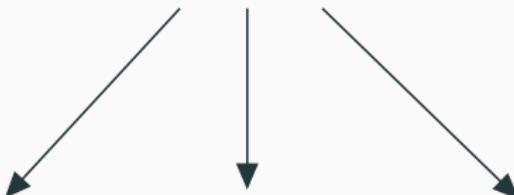
BIOINFORMATICS, 2018

TerryFoxLaboratory

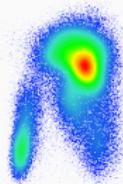
SFU



Efficient grouping methods for  
the annotation and sorting of single cells



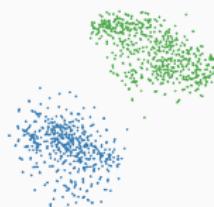
Flow cytometry



Metagenomics



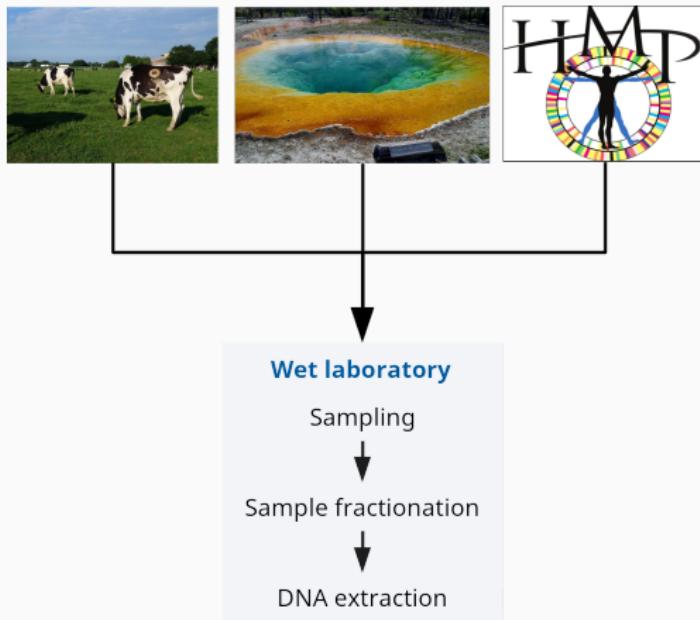
Single-cell genomics



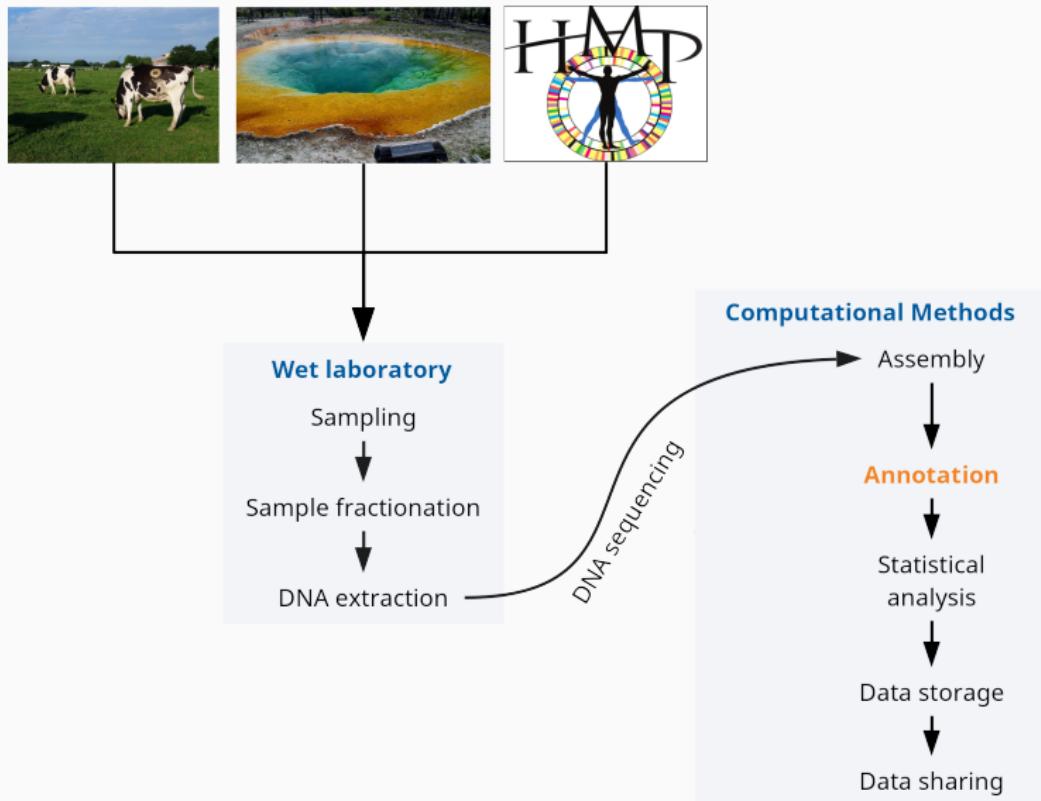
# METAGENOMICS IN A NUTSHELL



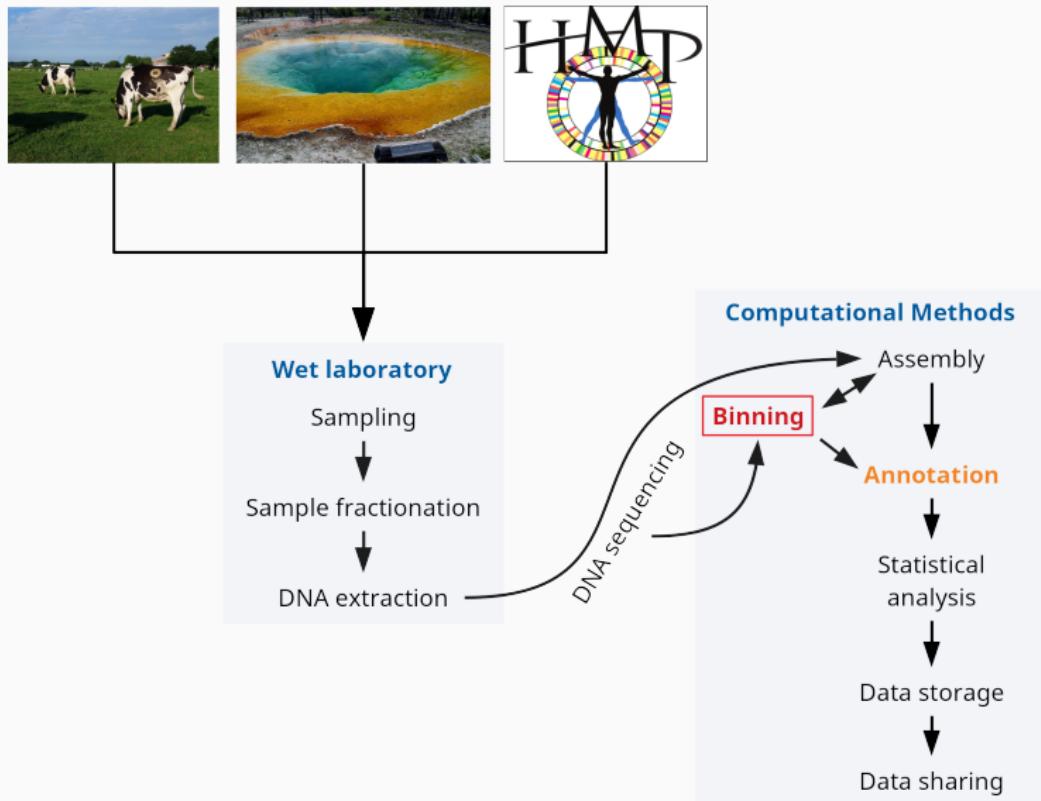
# METAGENOMICS IN A NUTSHELL



# METAGENOMICS IN A NUTSHELL



# METAGENOMICS IN A NUTSHELL



## PROBLEM STATEMENT

---

Binning:

What sequences belong to the same genome?

---

<sup>1</sup>Sedlar et al. (2017)

# PROBLEM STATEMENT

Binning:

What sequences belong to the same genome?

supervised



---

<sup>1</sup>Sedlar et al. (2017)

## PROBLEM STATEMENT

Binning:

What sequences belong to the same genome?

supervised



taxonomy-free<sup>1</sup>

abundance  
based

sequence  
composition  
based

---

<sup>1</sup>Sedlar et al. (2017)

## CONTRIBUTION

---

- Tools utilize sequence-composition based binning  
Laczny et al. (2014); Wu et al. (2015); Lin and Liao (2016), ...
- Integration of novel machine learning techniques
- Sensible w.r.t. involved techniques and parameters

## CONTRIBUTION

---

- Tools utilize sequence-composition based binning  
Laczny et al. (2014); Wu et al. (2015); Lin and Liao (2016), ...
  - Integration of novel machine learning techniques
  - Sensible w.r.t. involved techniques and parameters
- Thorough evaluation of pipeline ingredients, Lux et al. (2015b)
- Foundations for single-cell analysis

## Ingredients:

- Vectorial representation
- Dimensionality reduction
- Clustering

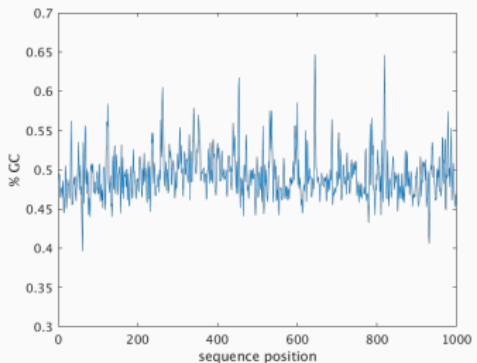
## Ingredients:

- Vectorial representation
- Dimensionality reduction
- Clustering

Idea: 1 cluster  $\leftrightarrow$  1 genome

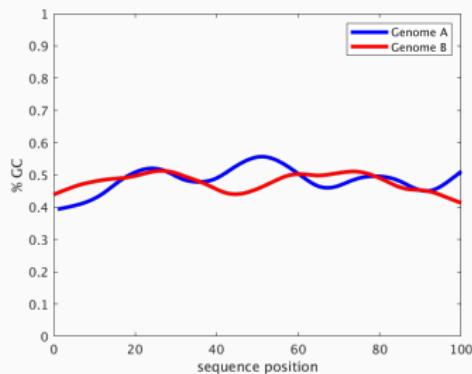
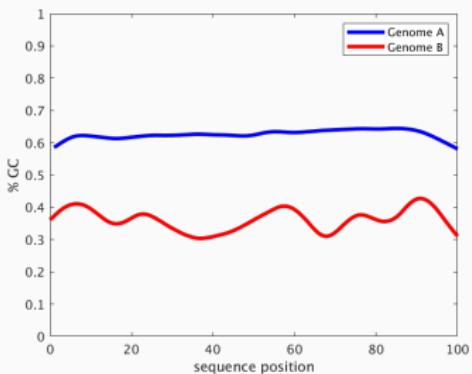
$\Rightarrow$  Resolve pronounced clusters

# GC-CONTENT AS A SEQUENCE SIGNATURE

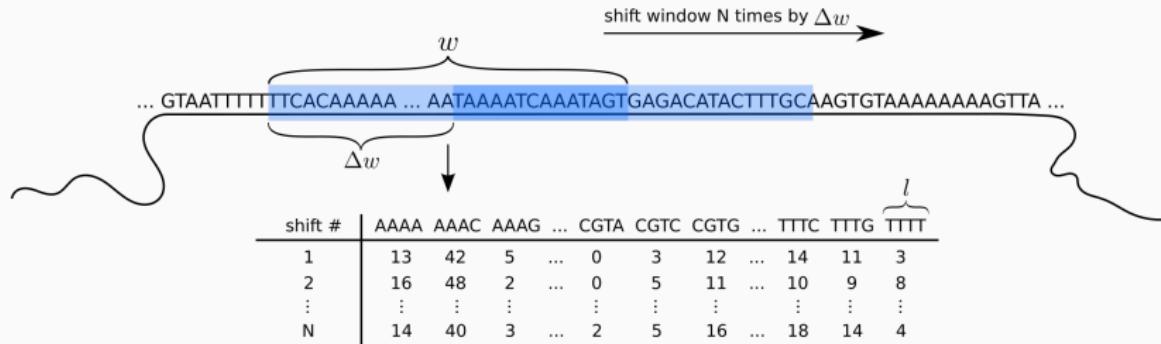


*E. coli* GC-content variation

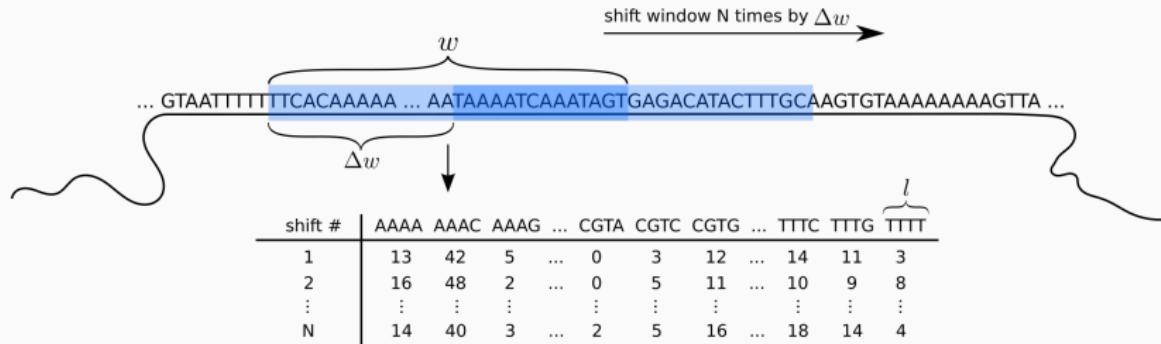
# GC-CONTENT AS A SEQUENCE SIGNATURE



# CALCULATION OF $l$ -MER FREQUENCIES



# CALCULATION OF $l$ -MER FREQUENCIES



## Observations:

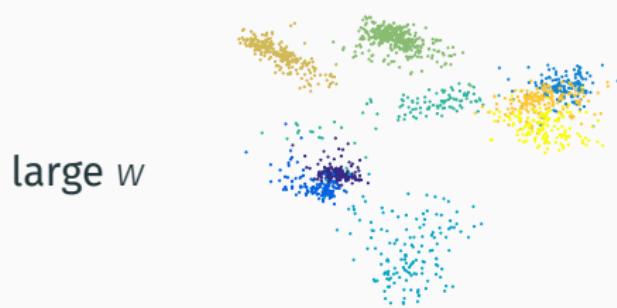
- $N$  data points with dimensionality  $4^l$
- Taking  $l = 1$  is equivalent to the GC-content
- Window parameters heavily influence representation

## EVALUATION QUESTIONS

---

- What window parameters to choose?
- How to reduce dimensionality?
- How to cluster?

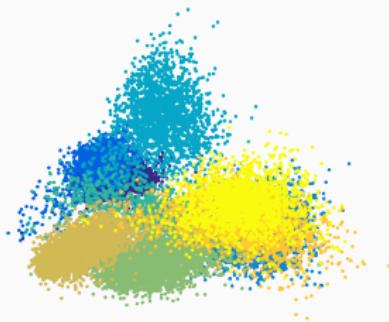
## DR & WINDOW PARAMETERS, QUALITATIVE



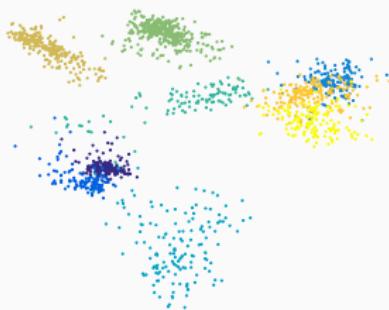
PCA

## DR & WINDOW PARAMETERS, QUALITATIVE

small  $w$



large  $w$

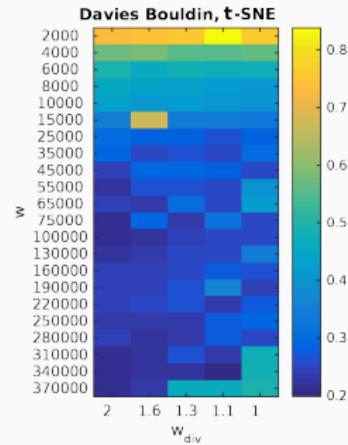
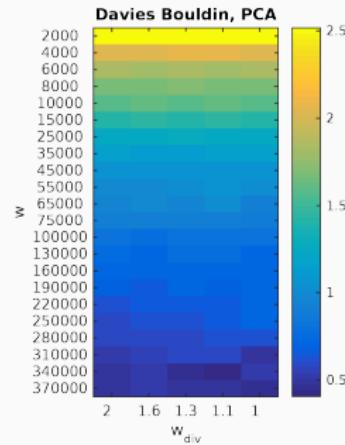
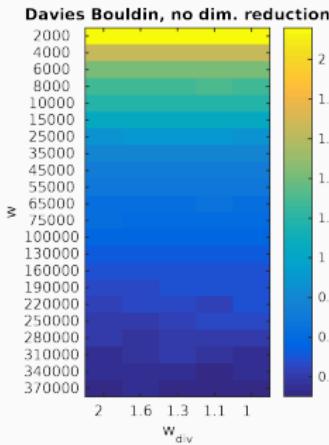


PCA



t-SNE

# DR & WINDOW PARAMETERS, QUANTITATIVE



→ Significant advantage of non-linear DR

Investigate suitable algorithms

How to determine number of clusters?

# CLUSTERING EVALUATION

THEO	DB		DUNN		XB		GAP		SS	
	K	J	K	J	K	J	K	J	K	K
$k_{opt} = 5$										
KM	5	<b>1.00</b>	5	<b>1.00</b>	3	0.56	55	0.10	<b>5</b>	
KM++	5	<b>1.00</b>	5	<b>1.00</b>	64	0.09	6			
NG	5	<b>1.00</b>	5	<b>1.00</b>	5	<b>1.00</b>	<b>5</b>			
HC	5	<b>1.00</b>	5	<b>1.00</b>	65	0.12	6			
GMM	5	<b>1.00</b>	5	<b>1.00</b>	7	0.80	<b>5</b>			
GMM++	5	<b>1.00</b>	5	<b>1.00</b>	-	-	12	0.47	7	
SC	5	<b>1.00</b>	5	<b>1.00</b>	-	-	-	-	20	
$K = 5, J = 1.00$										
DBS										
SC-EIG										
$k_{opt} = 10$										
KM	12	0.90	9	0.92	-	-	29	0.36	2	
KM++	<b>10</b>	<b>1.00</b>	10	<b>1.00</b>	10	0.00	32	0.34	<b>10</b>	
NG	9	0.93	10	1.00	11	1.00	9	0.93	2	
HC	<b>10</b>	<b>1.00</b>	10	<b>1.00</b>	26	0.46	14			
GMM	12	0.95	10	<b>1.00</b>	-	-	13	0.78	1	
GMM++	<b>10</b>	<b>1.00</b>	10	<b>1.00</b>	-	-	15	0.71	11	
SC	<b>10</b>	<b>1.00</b>	10	<b>1.00</b>	10	<b>1.00</b>	-	-	13	
$K = 10, J = 1.00$										
DBS										
SC-EIG										
$k_{opt} = 20$										
KM	23	0.79	8	0.36	-	-	51	0.42	2	
KM++	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	56	0.44	25	
NG	20	0.89	15	0.74	18	0.86	18	0.80	5	
HC	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	56	0.54	39	
GMM	28	0.81	5	0.22	-	-	13	0.55	1	
GMM++	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	<b>20</b>	<b>1.00</b>	40	0.58	23	
SC	33	0.55	<b>20</b>	<b>1.00</b>	7	0.11	-	-	2	
$K = 20, J = 1.00$										
DBS										
SC-EIG										
$k_{opt} = 30$										
KM	27	0.80	19	0.63	-	-	67	0.54	1	
KM++	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	68	0.50	39	
NG	26	0.86	34	0.94	32	0.94	24	0.83	3	
HC	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	49	0.79	62	
GMM	25	0.73	23	0.59	-	-	9	0.24	1	
GMM++	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	44	0.80	35	
SC	33	0.93	31	0.97	31	0.97	-	-	31	
$K = 30, J = 1.00$										
DBS										
SC-EIG										
$k_{opt} = 40$										
KM	39	0.66	21	0.49	-	-	68	0.65	1	
KM++	39	0.96	35	0.77	35	0.77	<b>40</b>	<b>1.00</b>	43	
NG	34	0.82	29	0.64	45	0.96	26	0.58	3	
HC	39	0.96	35	0.77	35	0.77	<b>40</b>	<b>1.00</b>	61	
GMM	44	0.62	9	0.21	-	-	6	0.15	3	
GMM++	39	0.96	35	0.77	35	0.77	<b>40</b>	<b>1.00</b>	44	
SC	35	0.56	3	0.03	2	0.28	-	-	41	
$K = 35, J = 0.77$										
DBS										
SC-EIG										
$k_{opt} = 50$										
KM	29	0.52	16	0.27	-	-	63	0.63	3	
KM++	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	54	
NG	35	0.63	51	0.83	43	0.74	24	0.45	11	
HC	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	70	
GMM	49	0.59	9	0.15	-	-	3	0.06	2	
GMM++	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	<b>50</b>	<b>1.00</b>	55	
SC	58	0.67	2	0.02	4	0.02	-	-	13	
$K = 36, J = 0.48$										
DBS										
SC-EIG										



NCBI	DB		DUNN		XB		GAP		SS	
	K	J	K	J	K	J	K	J	K	K
$k_{opt} = 5$										
KM	4	0.87	4	0.87	4	0.87	12	0.38	<b>5</b>	
KM++	<b>5</b>	<b>1.00</b>	4	0.87	4	0.87	11	0.41	7	
NG	<b>5</b>	<b>1.00</b>	4	0.87	4	0.87	10	0.45	8	
HC	<b>5</b>	<b>1.00</b>	4	0.87	4	0.87	6	0.80	6	
GMM	<b>5</b>	0.98	4	0.87	-	-	7	0.69	<b>5</b>	
GMM++	<b>5</b>	0.98	4	0.87	4	0.87	-	-	8	
SC	4	0.73	4	0.73	4	0.73	-	-	10	
$K = 4, J = 0.87$										
DBS										
SC-EIG										
$k_{opt} = 10$										
KM	<b>10</b>	<b>0.98</b>	6	0.63	9	0.94	68	0.13	4	
KM++	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	53	0.16	11	
NG	<b>10</b>	<b>0.98</b>	6	0.63	8	0.89	17	0.50	12	
HC	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	18	0.50	11	
GMM	9	0.79	8	0.71	-	-	13	0.89	9	
GMM++	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	<b>10</b>	<b>0.98</b>	50	0.20	11	
SC	2	0.13	2	0.13	2	0.13	-	-	12	
$K = 10, J = 0.98$										
DBS										
SC-EIG										
$k_{opt} = 20$										
KM	22	0.85	5	0.32	-	-	61	0.26	3	
KM++	22	0.73	16	0.91	16	0.90	63	0.24	<b>20</b>	
NG	<b>20</b>	0.74	4	0.27	13	0.76	30	0.53	18	
HC	<b>19</b>	<b>0.97</b>	17	0.90	18	0.96	39	0.44	19	
GMM	18	0.91	16	0.85	-	-	16	0.88	1	
GMM++	16	0.87	17	0.94	17	0.92	34	0.48	<b>20</b>	
SC	14	0.61	12	0.46	11	0.44	-	-	21	
$K = 18, J = 0.96$										
DBS										
SC-EIG										
$k_{opt} = 30$										
KM	<b>39</b>	0.75	17	0.53	-	-	70	0.53	4	
KM++	46	0.75	25	0.70	24	0.66	68	0.49	4	
NG	37	0.81	17	0.53	16	0.51	31	0.78	14	
HC	41	0.92	11	0.35	36	<b>0.95</b>	53	0.69	48	
GMM	29	0.78	17	0.49	4	0.12	17	0.45	1	
GMM++	43	0.81	29	0.80	4	0.12	-	-	1	
SC	31	0.6	4	0.04	7	0.06	-	-	1	
$K = 28, J = 0.91$										
DBS										
SC-EIG										
$k_{opt} = 40$										
KM	<b>40</b>	0.76								

## CLUSTERING EVALUATION – MAIN FINDINGS

---

- Majority of clustering methods work well on theoretical data
- Few methods deliver accurate results on simulated data
- Promising candidates:
  - Hierarchical Clustering
  - Connected Components Clustering
  - Davies-Bouldin index
- More complex metagenomes are more difficult

## CLUSTERING EVALUATION – MAIN FINDINGS

---

- Majority of clustering methods work well on theoretical data
- Few methods deliver accurate results on simulated data
- Promising candidates:
  - Hierarchical Clustering
  - Connected Components Clustering
  - Davies-Bouldin index
- More complex metagenomes are more difficult

## CLUSTERING EVALUATION – MAIN FINDINGS

---

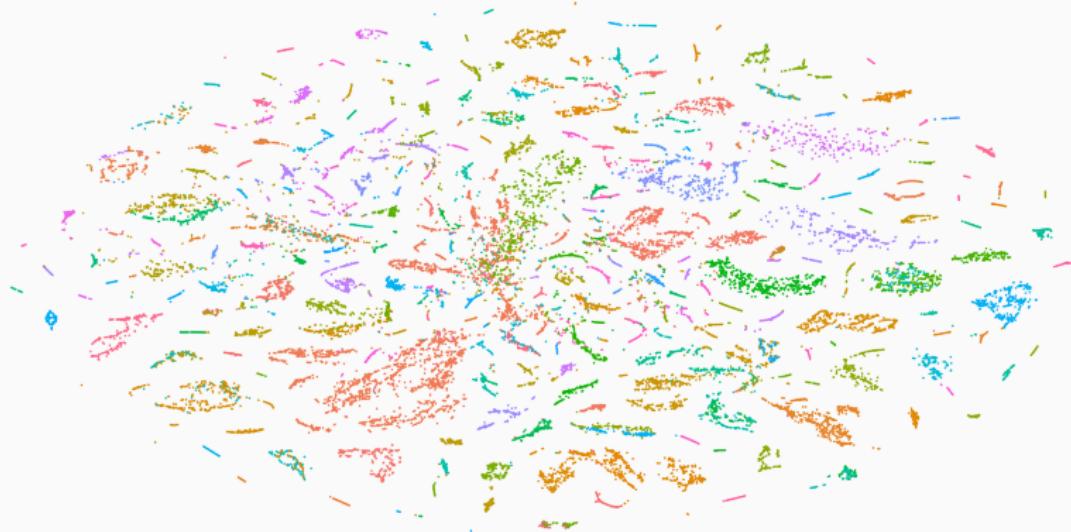
- Majority of clustering methods work well on theoretical data
- Few methods deliver accurate results on simulated data
- Promising candidates:
  - Hierarchical Clustering
  - Connected Components Clustering
  - Davies-Bouldin index
- More complex metagenomes are more difficult

## CLUSTERING EVALUATION – MAIN FINDINGS

---

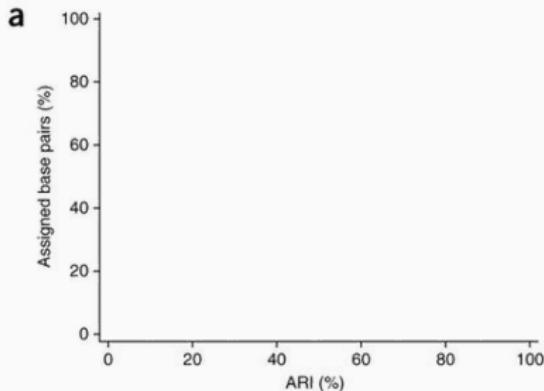
- Majority of clustering methods work well on theoretical data
- Few methods deliver accurate results on simulated data
- Promising candidates:
  - Hierarchical Clustering
  - Connected Components Clustering
  - Davies-Bouldin index
- More complex metagenomes are more difficult

## APPLICATION TO COMPLEX METAGENOMES



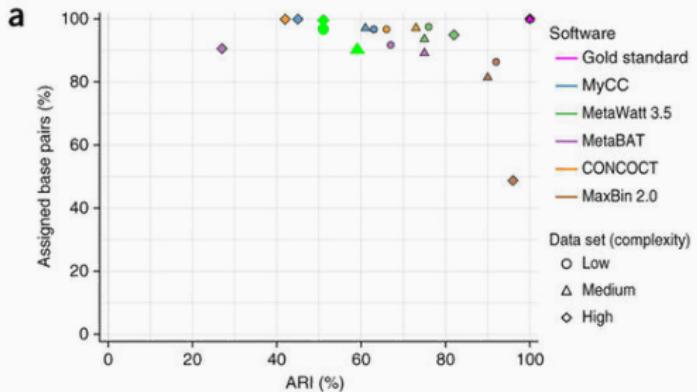
High complexity CAMI metagenome t-SNE representation,  
colored by gold standard

## APPLICATION TO COMPLEX METAGENOMES



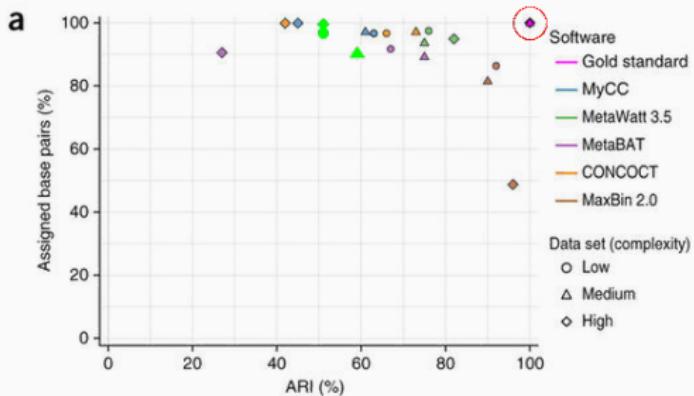
Comparison to other CAMI contestants

# APPLICATION TO COMPLEX METAGENOMES



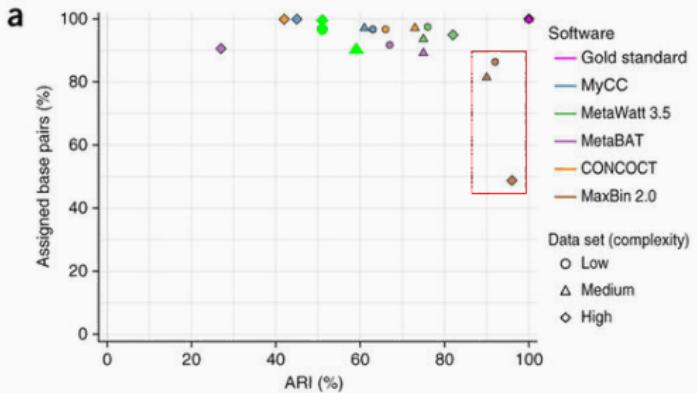
Comparison to other CAMI contestants

# APPLICATION TO COMPLEX METAGENOMES



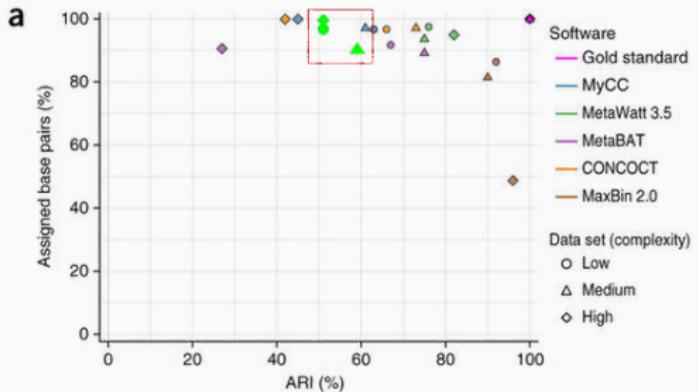
Comparison to other CAMI contestants

# APPLICATION TO COMPLEX METAGENOMES



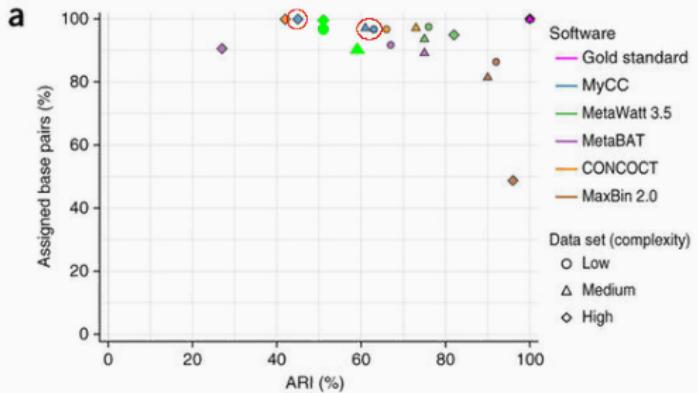
Comparison to other CAMI contestants

# APPLICATION TO COMPLEX METAGENOMES



Comparison to other CAMI contestants

# APPLICATION TO COMPLEX METAGENOMES



Comparison to other CAMI contestants

## GUIDELINES

---

- Choose large, half-overlapping windows for a fixed number of data points.
- Reduce dimensionality using cluster-focused t-SNE
- Hierarchical clustering + Davies-Bouldin index

# OUTLOOK

- Crucial: utilize auxiliary information
- Incorporate findings to extant tools
- Iterative clustering process
- Investigate recent DR alternatives (UMAP<sup>2</sup>)

---

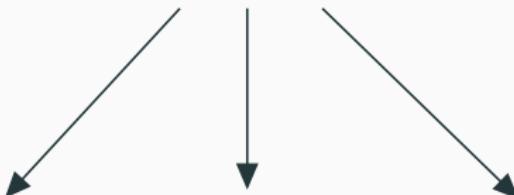
<sup>2</sup>McInnes and Healy (2018)

MARKUS LUX, ALEXANDER SCZYRBA, BARBARA HAMMER

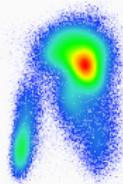
*AUTOMATIC DISCOVERY OF METAGENOMIC STRUCTURE*

2015 INTERNATIONAL JOINT CONFERENCE ON NEURAL NETWORKS (IJCNN)

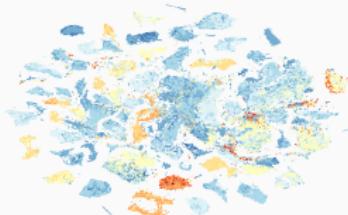
Efficient grouping methods for  
the annotation and sorting of single cells



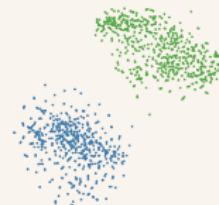
Flow cytometry



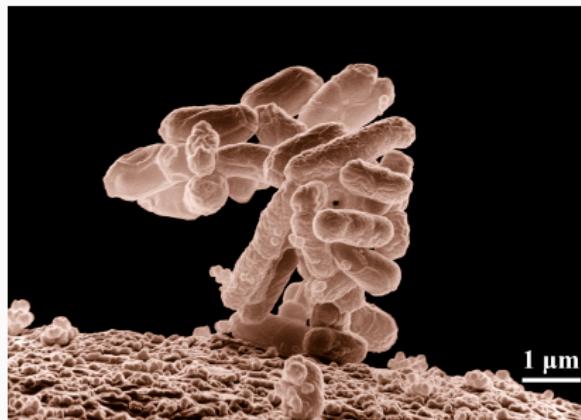
Metagenomics



Single-cell genomics

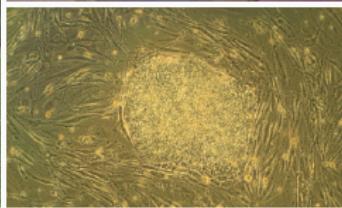
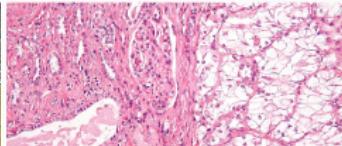


# SINGLE-CELL SEQUENCING



- Eliminates need for cultivation
- Amplification of one DNA molecule
- Focus on specifics of individual single cells

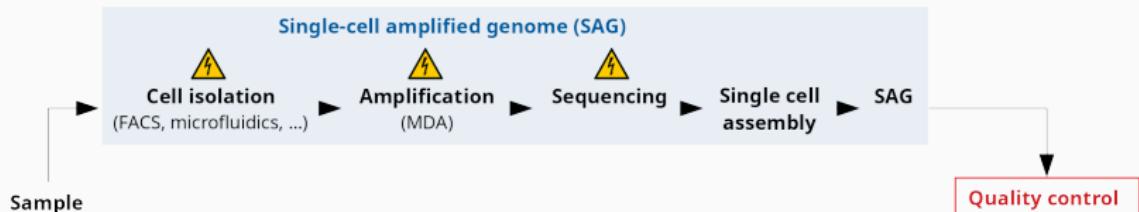
# SINGLE-CELL GENOMICS APPLICATIONS



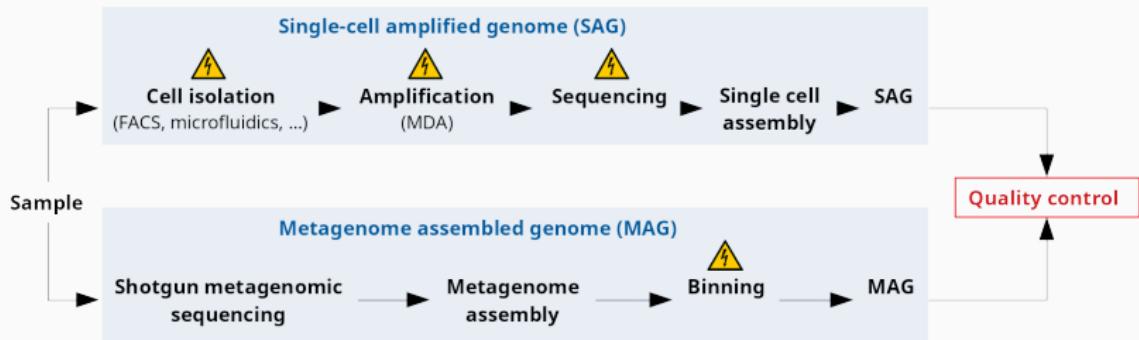
What if:

All these applications are based on wrong data?

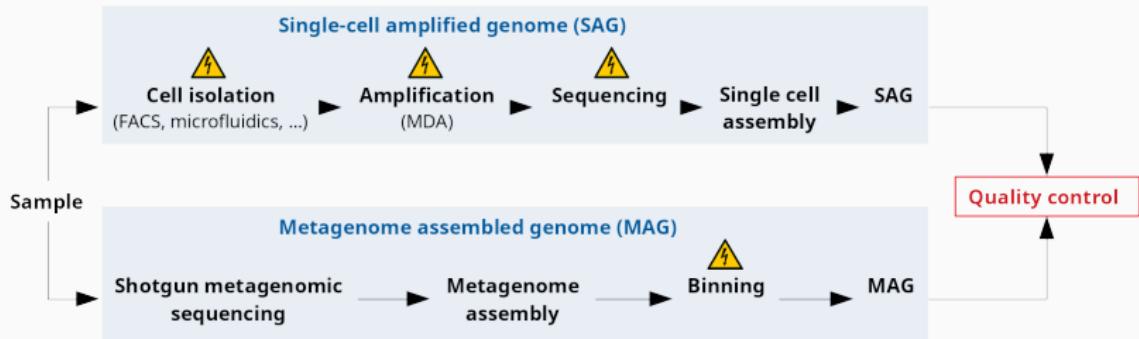
# PROBLEM: CONTAMINATED SAMPLES



# PROBLEM: CONTAMINATED SAMPLES



# PROBLEM: CONTAMINATED SAMPLES



- Sources of contamination:
  - In-sample
  - Unclean laboratory environment
  - Library preparation kits, reagents, ...
- Prevent introduction to public databases!

# QUALITY STANDARDS FOR REPORTING SAGs AND MAGs

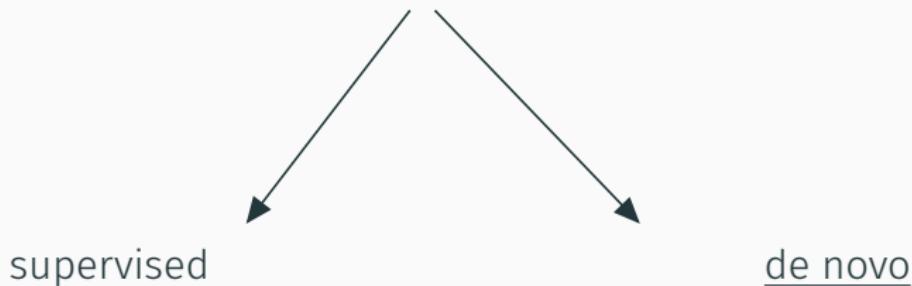
Contamination rates<sup>3</sup>:

- < 10% for low-quality draft genomes
- < 5% for high-quality draft genomes

---

<sup>3</sup>Bowers et al. (2017)

## Detecting contamination



ProDeGe, Tennesen et al. (2015)

# CONTRIBUTION

---

Interactive detection and  
cleansing without re-sequencing

# CONTRIBUTION

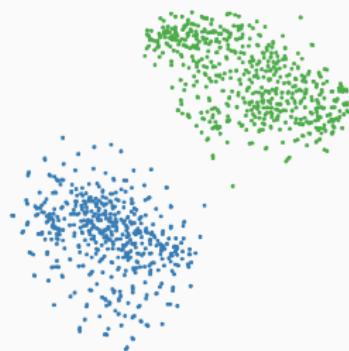
---

Interactive detection and  
cleansing without re-sequencing

acdc – automated contamination detection and confidence  
estimation for single-cell genome data

Lux et al. (2015a, 2016)

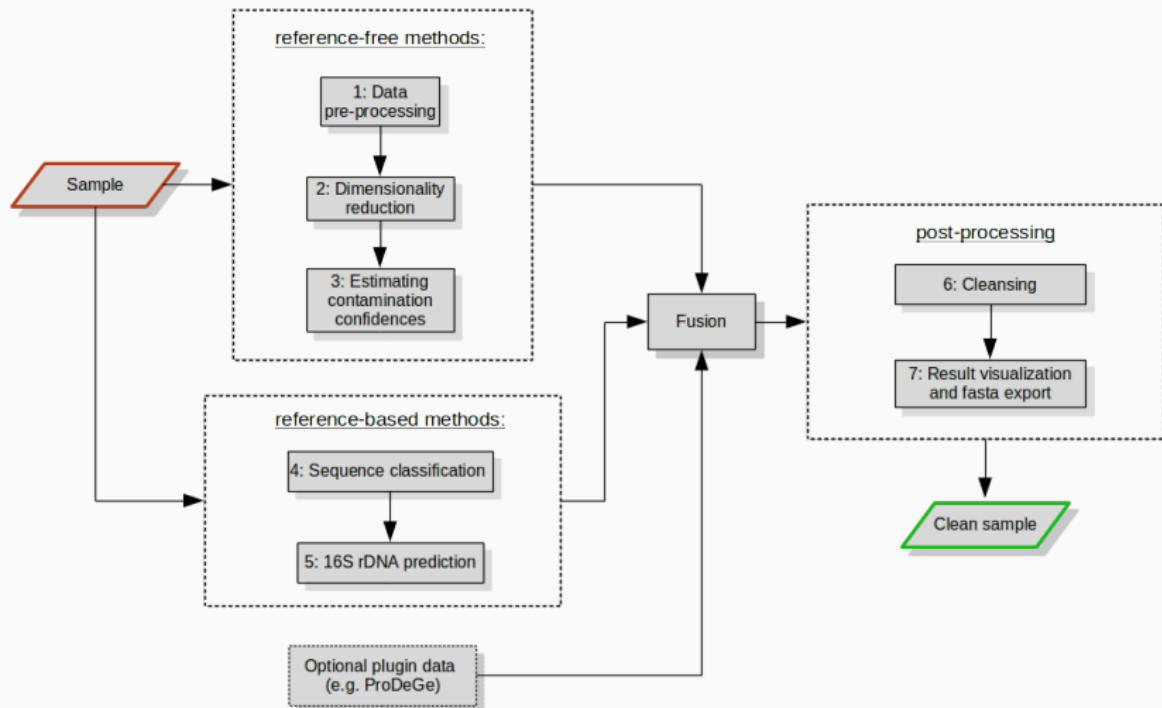
Unsupervised machine learning and clustering techniques.



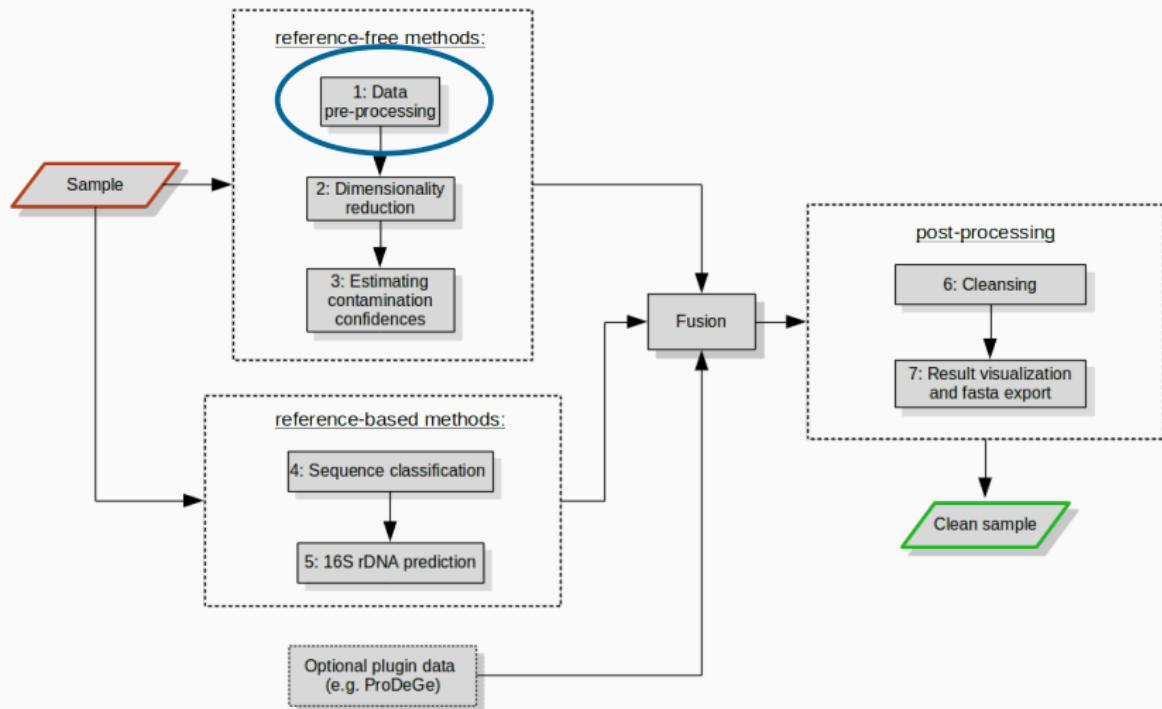
Idea: one cluster  $\Leftrightarrow$  one genome

→ detect number of clusters

# CONTAMINATION DETECTION AND CLEANSING



# CONTAMINATION DETECTION AND CLEANSING

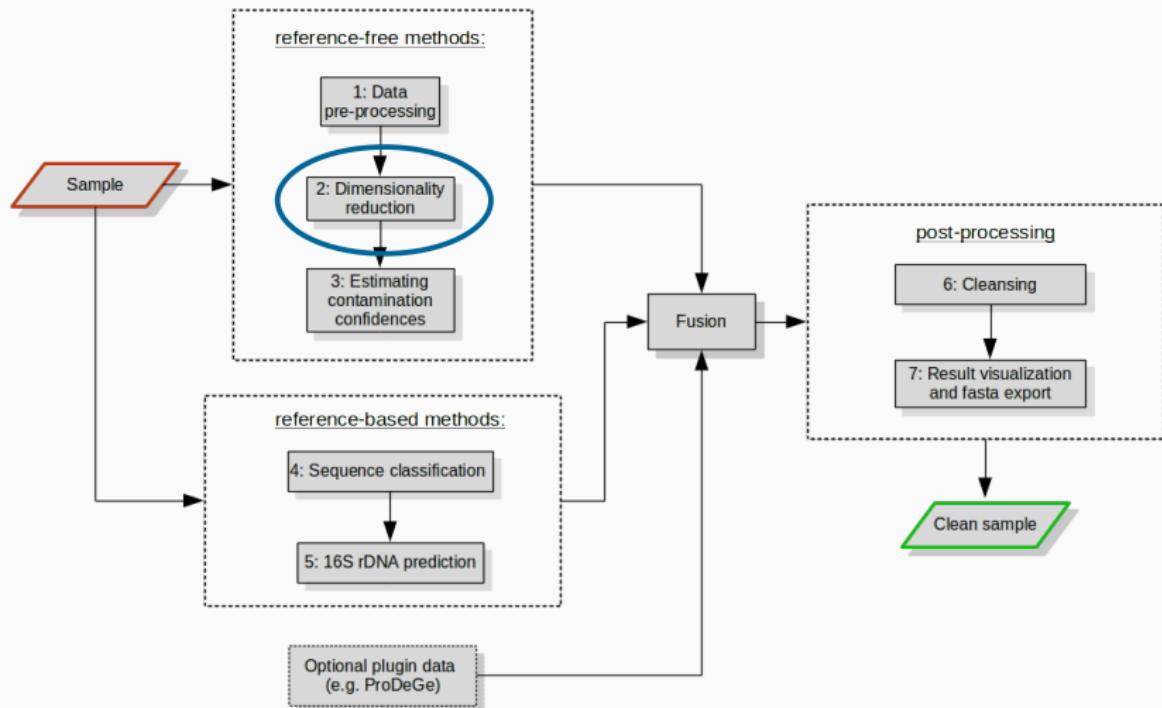


## WINDOW PARAMETERS

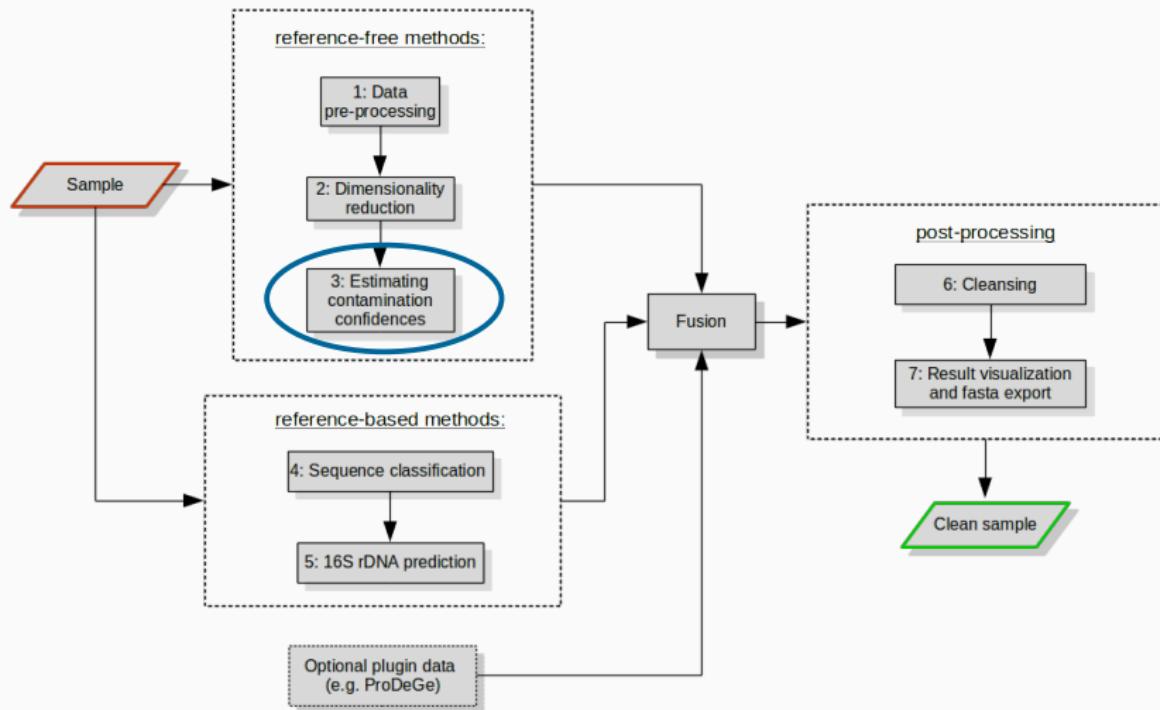


$$\Delta w = \lceil \frac{n_{bp}}{n_{target}} \rceil, w = 2 \cdot \Delta w$$

# CONTAMINATION DETECTION AND CLEANSING



# CONTAMINATION DETECTION AND CLEANSING



## CONTAMINATION YES/NO?

- Task: Determine whether number of clusters
  - $k = 1$  (clean)
  - $k > 1$  (contaminated)
- Many known algorithms not applicable<sup>4</sup>

---

<sup>4</sup>Lux et al. (2015a)

## CONTAMINATION YES/NO?

- Task: Determine whether number of clusters
  - $k = 1$  (clean)
  - $k > 1$  (contaminated)
- Many known algorithms not applicable<sup>4</sup>

### Contamination detection and confidence estimation:

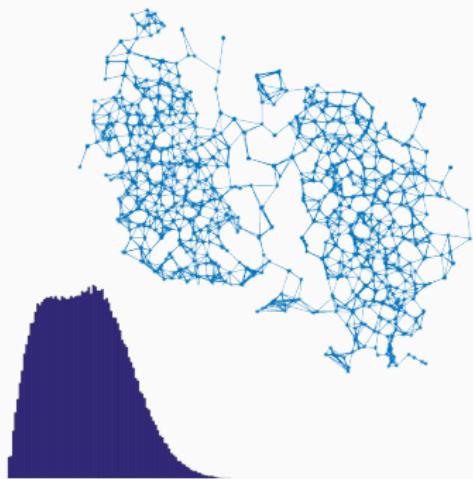
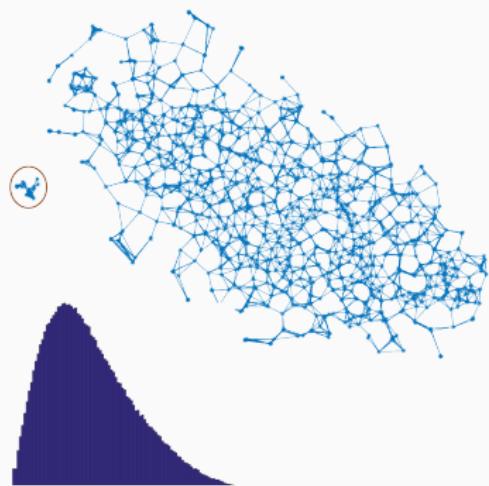
1. The dip-statistic test for multimodality
2. Connected components in a neighborhood graph
3. Estimate confidences using bootstrapping

---

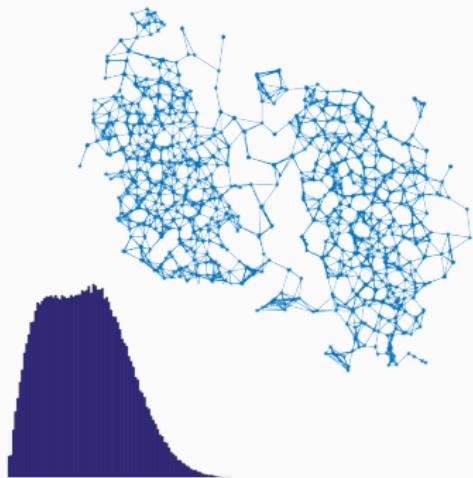
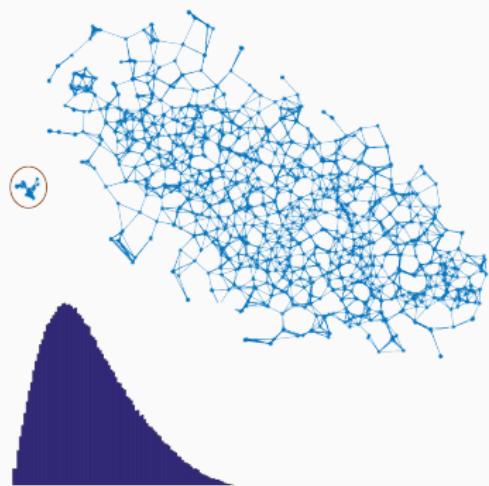
<sup>4</sup>Lux et al. (2015a)

CONTAMINATION YES/NO?

---



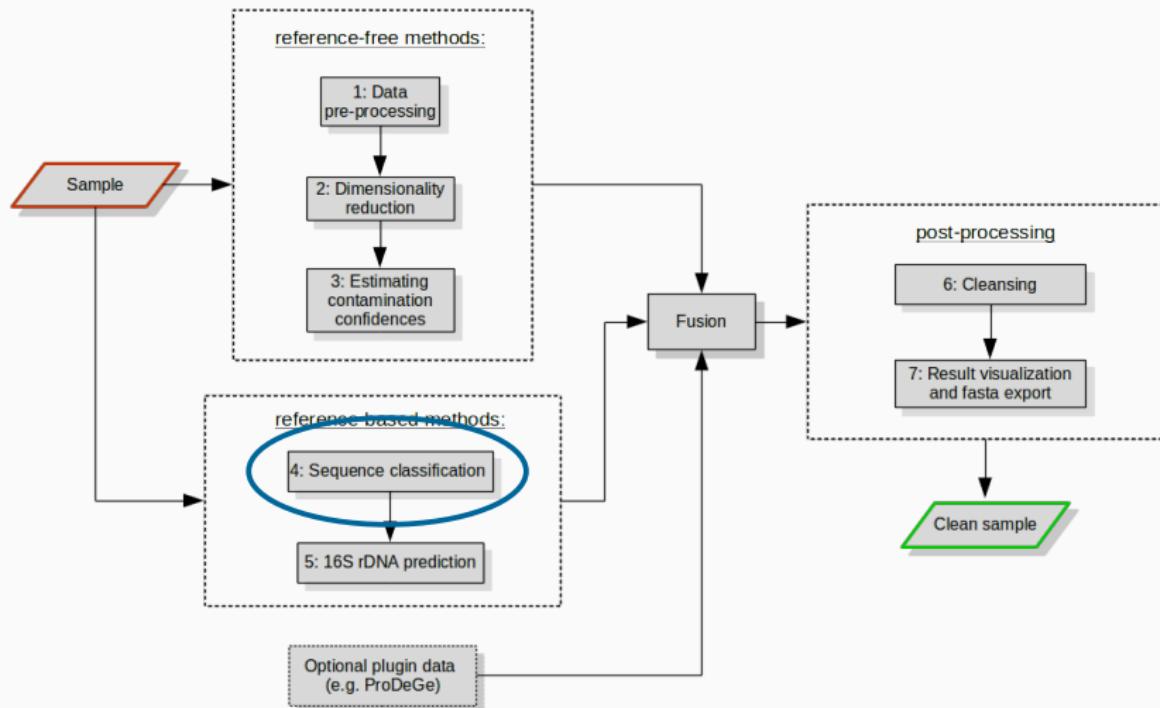
# CONTAMINATION YES/NO?



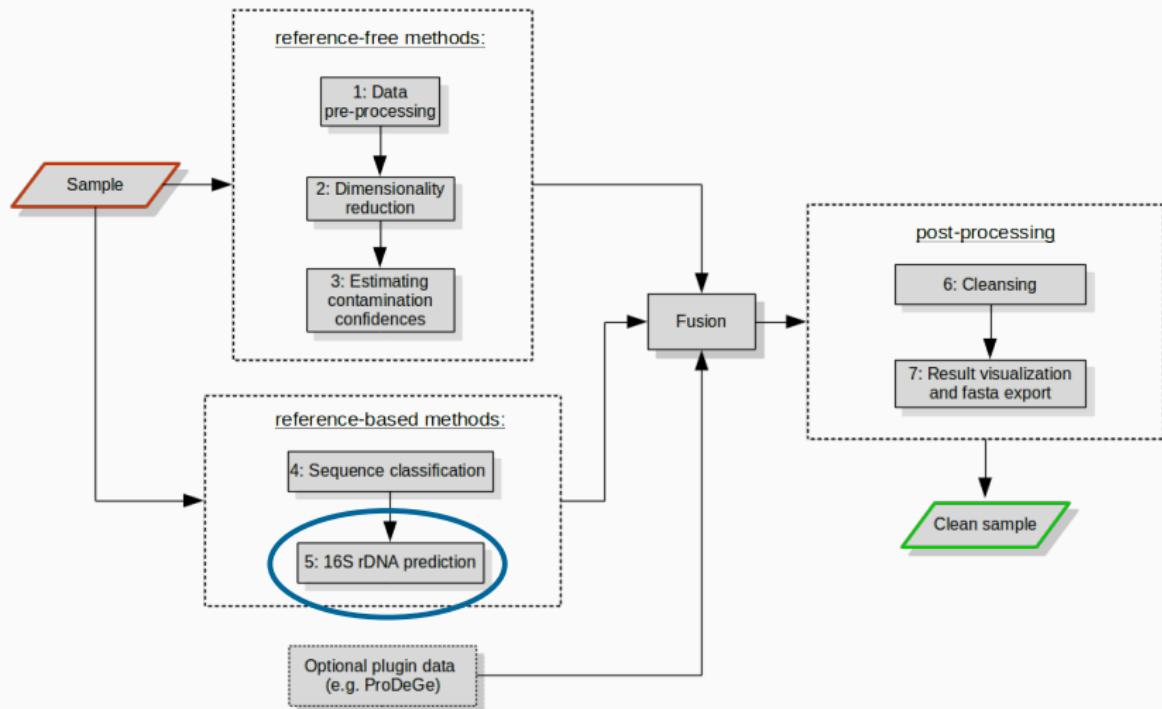
Confidence value:

$$\nu_{dip|cc} = \frac{1}{B} \sum_{b=1}^B \mathbf{1}_{k_{dip} + k_{cc} > 2}$$

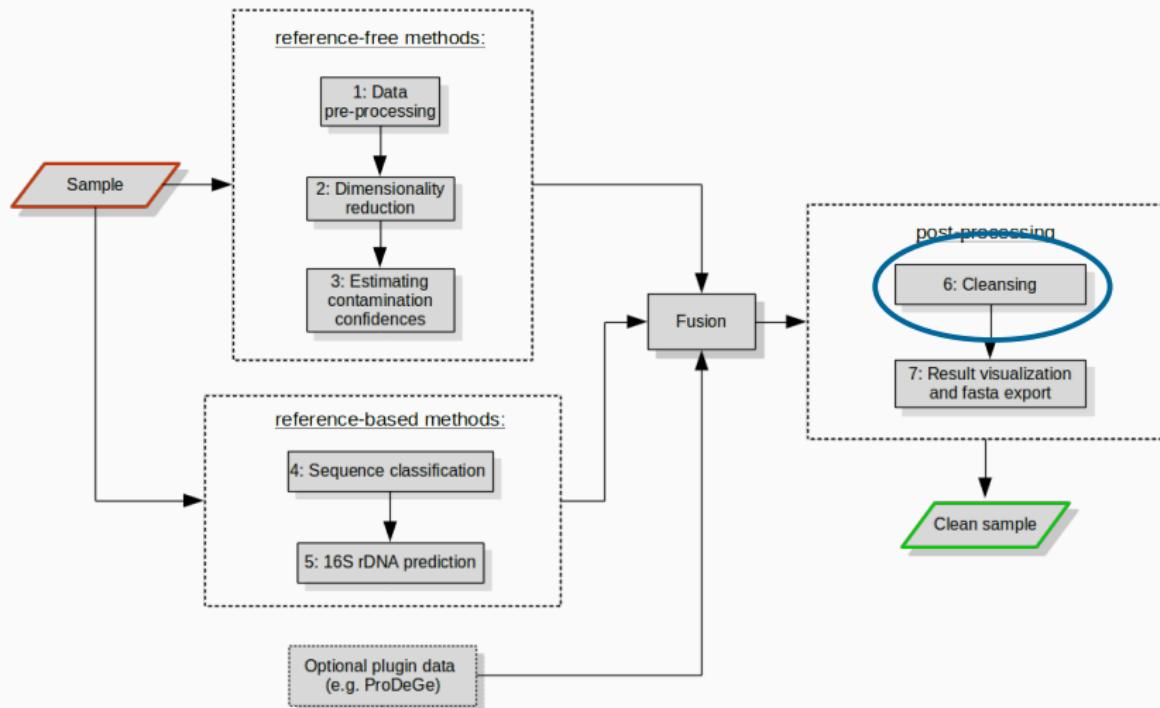
# CONTAMINATION DETECTION AND CLEANSING

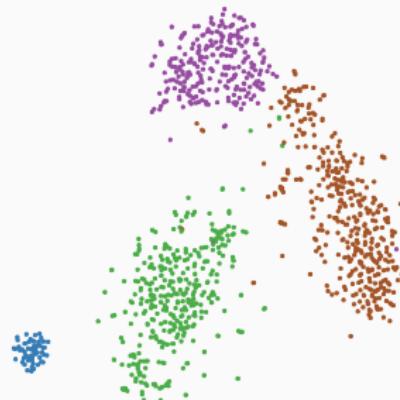


# CONTAMINATION DETECTION AND CLEANSING



# CONTAMINATION DETECTION AND CLEANSING





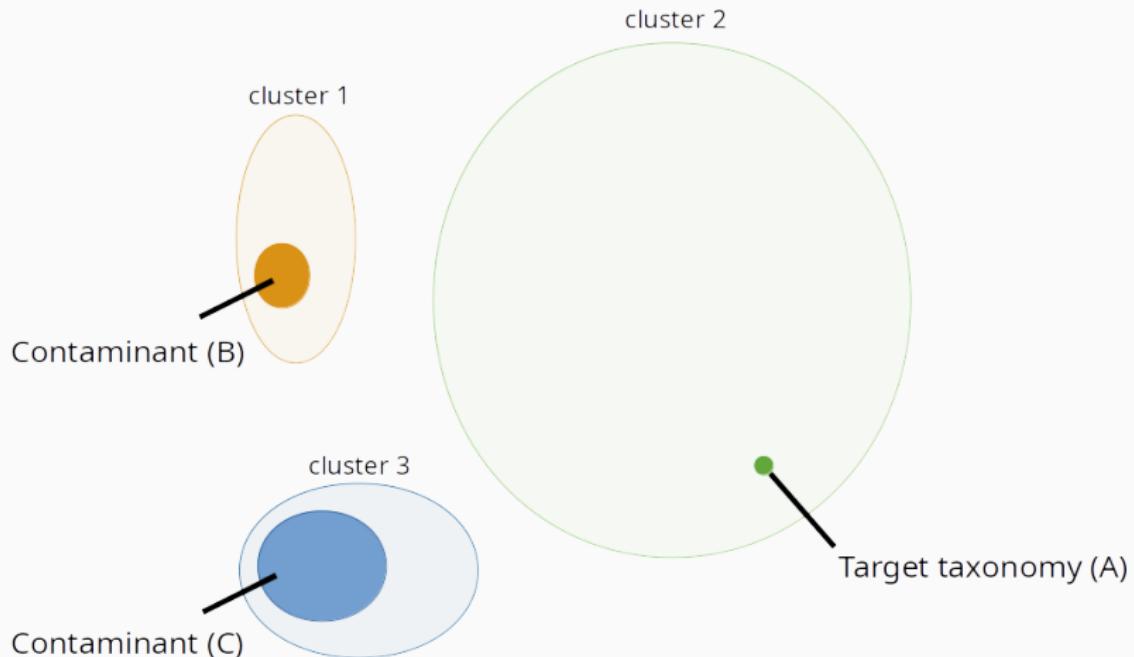
Find optimal clustering for a given range of  $k$

→ export individual clusters

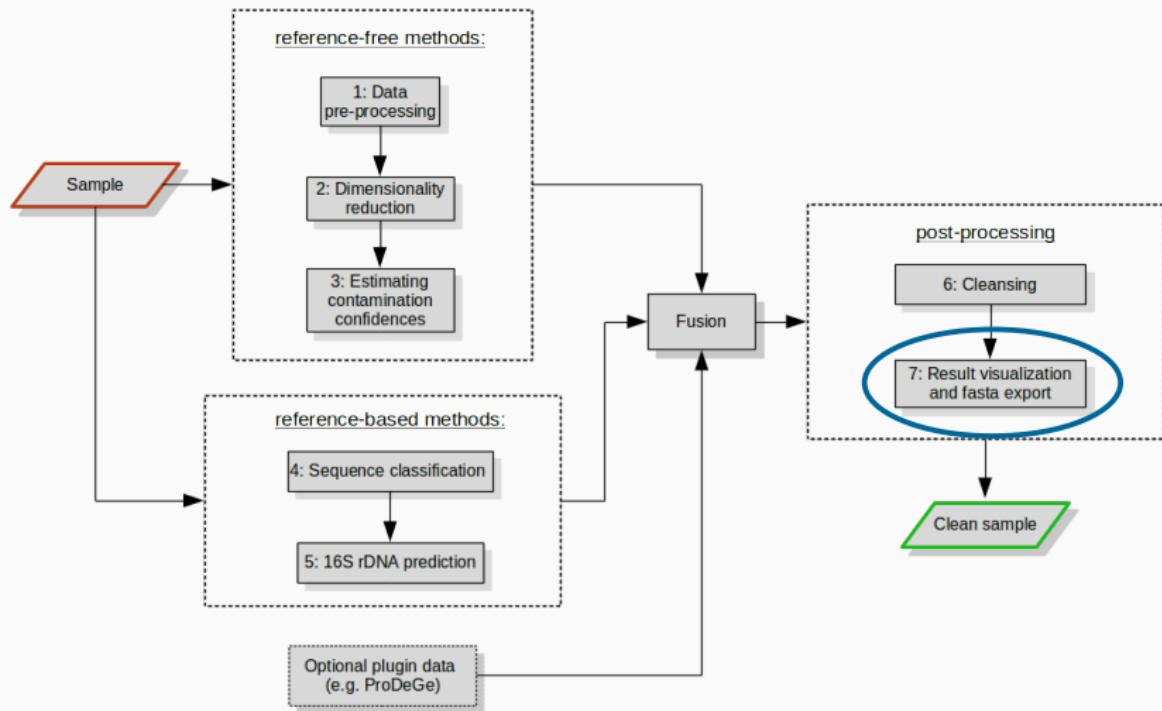
# INTEGRATION OF EXTERNAL KNOWLEDGE: TAXONOMY INFERENCE

---

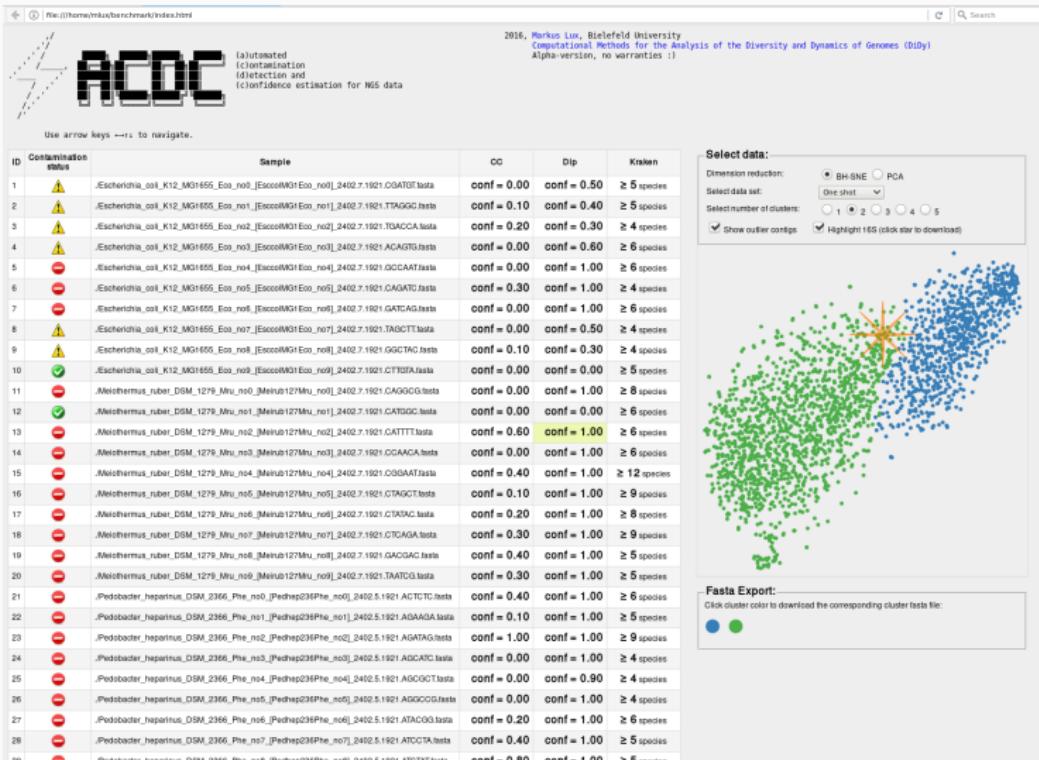
# INTEGRATION OF EXTERNAL KNOWLEDGE: TAXONOMY INFERENCE



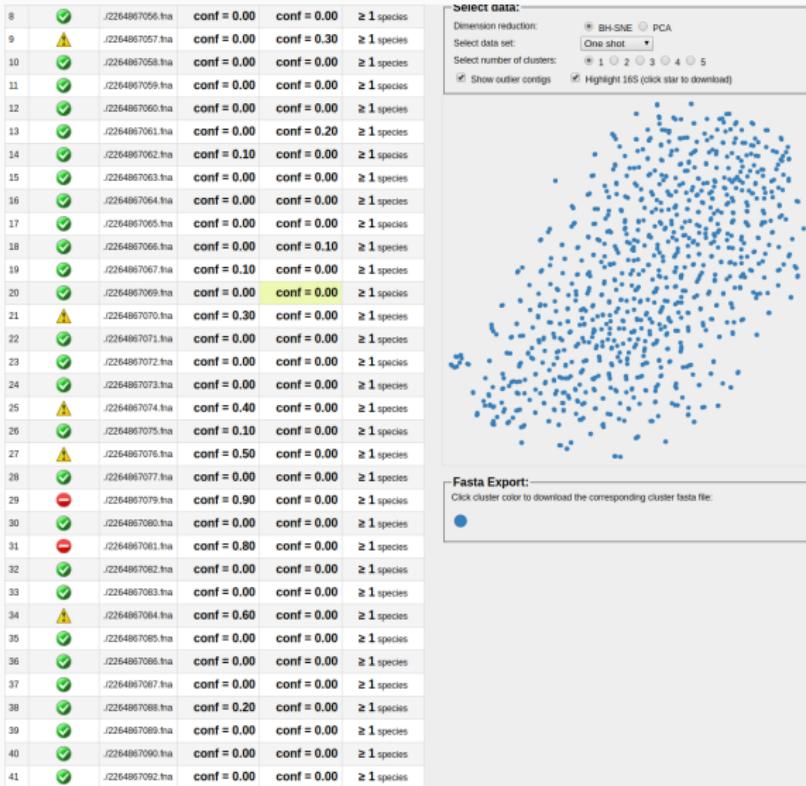
# CONTAMINATION DETECTION AND CLEANSING



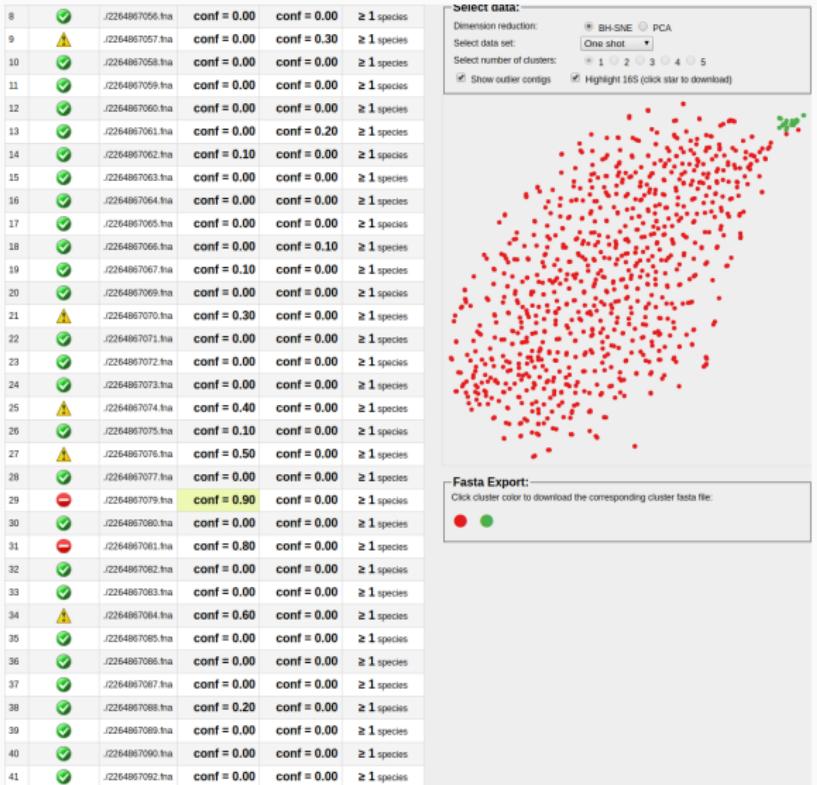
# EXAMPLE 1



## EXAMPLE 2



# CLEAN SAMPLES ACTUALLY CONTAMINATED



## RESULTS

---

- Correctly detected contamination state on large set of SAGs
- Contaminants of same genus
- Removed contaminants with high precision and recall
- Good prediction of pure and impure CAMI bins
- Quadratic complexity but low practical requirements

## RESULTS

---

- Correctly detected contamination state on large set of SAGs
- Contaminants of same genus
- Removed contaminants with high precision and recall
- Good prediction of pure and impure CAMI bins
- Quadratic complexity but low practical requirements

## RESULTS

---

- Correctly detected contamination state on large set of SAGs
- Contaminants of same genus
- Removed contaminants with high precision and recall
- Good prediction of pure and impure CAMI bins
- Quadratic complexity but low practical requirements

## RESULTS

---

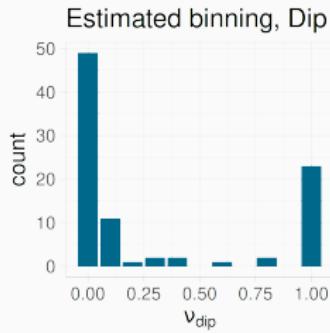
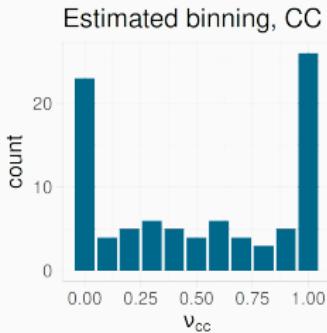
- Correctly detected contamination state on large set of SAGs
- Contaminants of same genus
- Removed contaminants with high precision and recall
- Good prediction of pure and impure CAMI bins
- Quadratic complexity but low practical requirements

## RESULTS

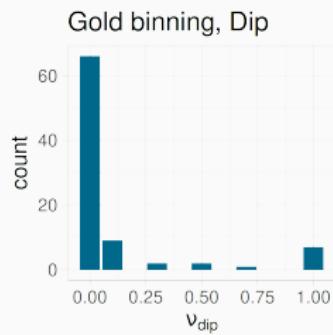
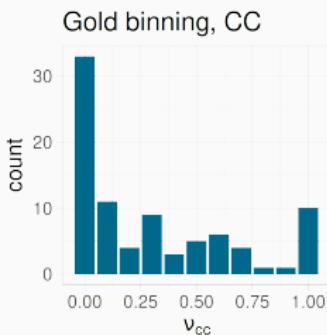
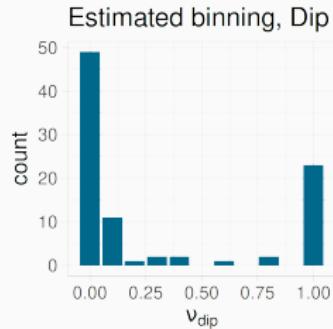
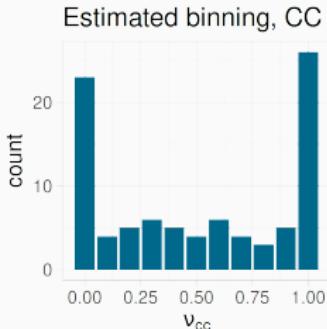
---

- Correctly detected contamination state on large set of SAGs
- Contaminants of same genus
- Removed contaminants with high precision and recall
- Good prediction of pure and impure CAMI bins
- Quadratic complexity but low practical requirements

# EARLY RESULTS ON MAG DATA



# EARLY RESULTS ON MAG DATA



# OUTLOOK

---

- Collaboration with Joint Genome Institute (JGI), USA
- Research other types of data representation to overcome species-rank limitation
- Further investigation of MAG capabilities
- Investigate horizontal gene transfer detection

## OUTLOOK

---

- Collaboration with Joint Genome Institute (JGI), USA
- Research other types of data representation to overcome species-rank limitation
- Further investigation of MAG capabilities
- Investigate horizontal gene transfer detection

## OUTLOOK

---

- Collaboration with Joint Genome Institute (JGI), USA
- Research other types of data representation to overcome species-rank limitation
- Further investigation of MAG capabilities
- Investigate horizontal gene transfer detection

## OUTLOOK

---

- Collaboration with Joint Genome Institute (JGI), USA
- Research other types of data representation to overcome species-rank limitation
- Further investigation of MAG capabilities
- Investigate horizontal gene transfer detection

MARKUS LUX, JAN KRÜGER, CHRISTIAN RINKE, IRENA MAUS, ANDREAS SCHLUETER, TANJA WOYKE, ALEXANDER SCZYRBA, BARBARA HAMMER

*ACDC – AUTOMATED CONTAMINATION DETECTION AND CONFIDENCE ESTIMATION FOR SINGLE-CELL GENOME DATA*

BMC BIOINFORMATICS, 2016

---

MARKUS LUX, BARBARA HAMMER, ALEXANDER SCZYRBA

*AUTOMATED CONTAMINATION DETECTION IN SINGLE-CELL SEQUENCING*

BIORxIV, 2015

---

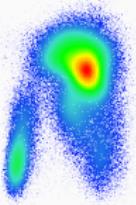
*POSTER @ THIRD MICROBIAL SINGLE CELL GENOMICS WORKSHOP*

BIGELOW LABORATORY FOR OCEAN SCIENCES, 2015

---



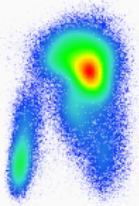
## OVERALL SUMMARY



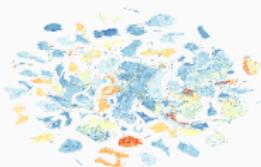
- Semi-supervised approach for flow cytometry gating
- Gate thousands of samples from one example, outperforming state-of-the-art methods

Lux et al. (2018)

## OVERALL SUMMARY

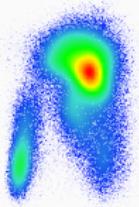


- Semi-supervised approach for flow cytometry gating
  - Gate thousands of samples from one example, outperforming state-of-the-art methods
- Lux et al. (2018)
- 

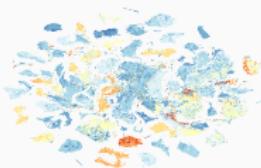


- Evaluation of taxonomy-independent, sequence-composition based binning
  - Guidelines for data representation, dimensionality reduction and clustering
- Lux et al. (2015b)

# OVERALL SUMMARY



- Semi-supervised approach for flow cytometry gating
  - Gate thousands of samples from one example, outperforming state-of-the-art methods
- Lux et al. (2018)
- 



- Evaluation of taxonomy-independent, sequence-composition based binning
  - Guidelines for data representation, dimensionality reduction and clustering
- Lux et al. (2015b)
- 



- Automated contamination detection for single-cell genome data
- Currently adapted for production by leading single-cell institute

Lux et al. (2015a, 2016)

# REFERENCES

- Nima Aghaeepour, Radina Nikolic, Holger H Hoos, and Ryan R Brinkman. Rapid cell population identification in flow cytometry data. *Cytometry Part A*, 79(1):6–13, 2011.
- Robert M Bowers, Nikos C Kyrides, Ramunas Stepanauskas, Miranda Harmon-Smith, Devin Doud, TBK Reddy, Frederik Schulz, Jessica Jarett, Adam R Rivers, Emiley A Eloe-Fadrosh, et al. Minimum information about a single amplified genome (misag) and a metagenome-assembled genome (mimag) of bacteria and archaea. *Nature biotechnology*, 35(8), 2017.
- Greg Finak, Jacob Frelinger, Wenxin Jiang, Evan W Newell, John Ramey, Mark M Davis, Spyros A Kalams, Stephen C De Rosa, and Raphael Gottardo. Opencyto: an open source infrastructure for scalable, robust, reproducible, and automated, end-to-end flow cytometry data analysis. *PLoS computational biology*, 10(8):e1003806, 2014.
- Greg Finak, Marc Langweiler, Maria Jaimes, Mehrnoush Malek, Jafar Taghiyar, Yael Korin, Khadir Raddassi, Lesley Devine, Gerlinde Obermoser, Marcin L Pekalski, et al. Standardizing flow cytometry immunophenotyping analysis from the human immunophenotyping consortium. *Scientific reports*, 6, 2016.
- Yongchao Ge and Stuart C Sealfon. flowpeaks: a fast unsupervised clustering for flow cytometry data via k-means and density peak finding. *Bioinformatics*, 28(15):2052–2058, 2012.
- Cedric C Laczny, Nicolás Pinel, Nikos Vlassis, and Paul Wilmes. Alignment-free visualization of metagenomic data by nonlinear dimension reduction. *Scientific reports*, 4, 2014.
- Huamin Li, Uri Shaham, Kelly P. Stanton, Yi Yao, Ruth Montgomery, and Yuval Kluger. Gating mass cytometry data by deep learning. *Bioinformatics*, 2017.
- Hsin-Hung Lin and Yu-Chieh Liao. Accurate binning of metagenomic contigs via automated clustering sequences using information of genomic signatures and marker genes. *Scientific reports*, 6:24175, 2016.

- Markus Lux, Barbara Hammer, and Alexander Sczyrba. Automated contamination detection in single-cell sequencing. *bioRxiv*, page 020859, 2015a.
- Markus Lux, Alexander Sczyrba, and Barbara Hammer. Automatic discovery of metagenomic structure. In *Neural Networks (IJCNN), 2015 International Joint Conference on*, pages 1–8. IEEE, 2015b.
- Markus Lux, Jan Krüger, Christian Rinke, Irena Maus, Andreas Schlüter, Tanja Woyke, Alexander Sczyrba, and Barbara Hammer. acdc—automated contamination detection and confidence estimation for single-cell genome data. *BMC bioinformatics*, 17(1):543, 2016.
- Markus Lux, Ryan Remy Brinkman, Cedric Chauve, Adam Laing, Anna Lorenc, Lucie Abeler-Dörner, Barbara Hammer, and Jonathan Wren. flowlearn: Fast and precise identification and quality checking of cell populations in flow cytometry. *Bioinformatics*, 1:9, 2018.
- Florian Mair, Felix J Hartmann, Dunja Mrdjen, Vinko Tosevski, Carsten Krieg, and Burkhard Becher. The end of gating? an introduction to automated analysis of high dimensional cytometry data. *European journal of immunology*, 46(1):34–43, 2016.
- Mehrnoosh Malek, Mohammad Jafar Taghiyar, Lauren Chong, Greg Finak, Raphael Gottardo, and Ryan R Brinkman. flowdensity: reproducing manual gating of flow cytometry data by automated density-based cell population identification. *Bioinformatics*, 31(4):606–607, 2015.
- Leland McInnes and John Healy. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Peng Qiu, Erin F Simonds, Sean C Bendall, Kenneth D Gibbs Jr, Robert V Bruggner, Michael D Linderman, Karen Sachs, Garry P Nolan, and Sylvia K Plevritis. Extracting a cellular hierarchy from high-dimensional cytometry data with spade. *Nature biotechnology*, 29(10):886–891, 2011.
- Karel Sedlar, Kristyna Kupkova, and Ivo Provaznik. Bioinformatics strategies for taxonomy independent binning and visualization of sequences in shotgun metagenomics. *Computational and structural biotechnology journal*, 15: 48–55, 2017.

- Kristin Tennessen, Evan Andersen, Scott Clingenpeel, Christian Rinke, Derek S Lundberg, James Han, Jeff L Dangl, Natalia Ivanova, Tanja Woyke, Nikos Kyrpides, et al. Prodege: a computational protocol for fully automated decontamination of genomes. *The ISME journal*, 2015.
- Sofie Van Gassen, Britt Callebaut, Mary J Van Helden, Bart N Lambrecht, Piet Demeester, Tom Dhaene, and Yvan Saeys. Flowsom: Using self-organizing maps for visualization and interpretation of cytometry data. *Cytometry Part A*, 87(7):636–645, 2015.
- Yu-Wei Wu, Blake A Simmons, and Steven W Singer. Maxbin 2.0: an automated binning algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics*, 32(4):605–607, 2015.