## rapids-singlecell

Getting your analysis done faster







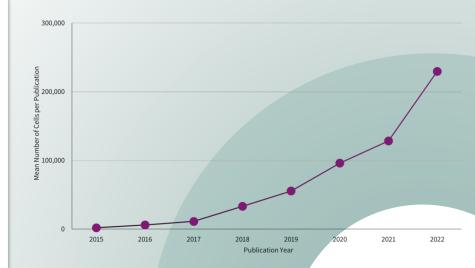




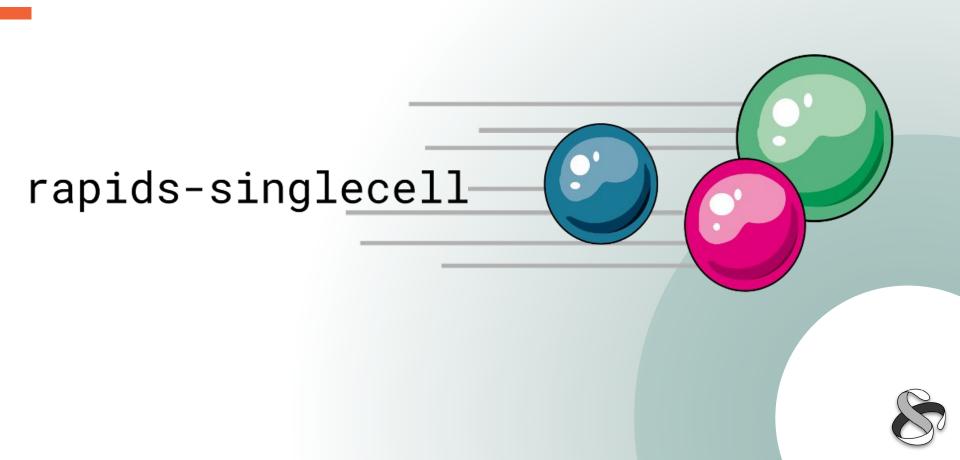


#### What's the need for another Package?

- Python is slow
- Python doesn't scale well
- People want to not do coffee breaks all day
- Maybe you want to work with collaborators in real time
- Singlecell analysis need a lot of iterations







#### What is rapids-singlecell?

- GPU accelerated singlecell analysis
- Near drop-in replacement for scanpy
- more than 30x Speedup over the CPU
- fully compatible with the scverse ecosystem



#### **Tech Stack RSC**



## CuPy

## RAP)DS



#### How do I get RSC?

- PyPi
- Conda/YAML
- Docker
- Pitfalls:
  - Wrong CUDA Stack
  - Issues with CuPy
- Solution:
  - Talk to your IT
  - Know you CUDA Version
  - Use CuPy Cuda specific wheels

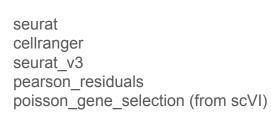


## What's Supported?



#### scanpy.pp

- filter cells / filter genes
- calculate\_qc\_metrics
- normalize\_total
- log1p
- highly\_variable\_genes
- regress\_out
- scale
- PCA
- Neighbors
- Harmony







#### scanpy.tl

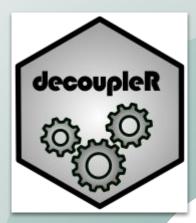
- Leiden
- Louvain
- TSNE
- UMAP
- Draw Graph
- Kernel Density
- Diffusion Maps
- Score Cells
- Rank Gene Groups with logreg





#### decoupler

- different statistical methods to extract biological activities
- run\_mlm
- run\_wsum
- run\_ulm
- run\_aucell





## squidpy

squidpy

- Spatial Autocorrelation
  - Moran's I
  - o Geary's C
- Ligrec (receptor-ligand analysis)



#### Get

#### **Utilities for AnnData Handling**

- aggregate
  - scanpy.get.aggregate
  - 'count\_nonzero', 'mean', 'sum', 'var'
  - o output dense and sparse
- anndata\_to\_GPU & anndata\_to\_CPU
- X\_to\_GPU & X\_to\_CPU



## Who uses/used rapids-singlecell?



#### **Memory Management**

Performance Tuning with RMM

- Memory management has an impact on performance
- Managed Memory:
  - Oversubscription
  - o more Data
  - less performant
- Pool Allocation:
  - Maximal Performance
  - more memory usage.



#### Limitations

- Not every function is ported to GPU yet
- Int32 indptr 2^31-1 nnz
- GPUs might be slower for small Datasets
- Memory Allocators
- Rapids Managed unified Memory doesn't allow for limitless oversubscription of VRAM



#### Outlook

Multi GPU with Dask

- Normalize
- Log1p
- Basic qc
- Scale
- PCA (dense from CUML +Sparse)
- HVG (Cellranger, Seurat)







Documentation



Github



#### **Contributions are Welcome**

- Tell me what's missing
- General improvements
- Start contributing yourself



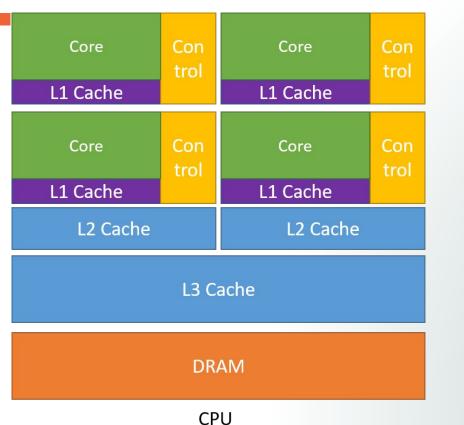
#### **Acknowledgments**

- scverse
  - Isaac Virshup
  - Lukas Heumos
  - Philipp Angerer
  - Ilan Gold
- Nvidia
  - Taurean Dyer
  - Corey Nolet
- CZI
- Fabian Theis
- The Users



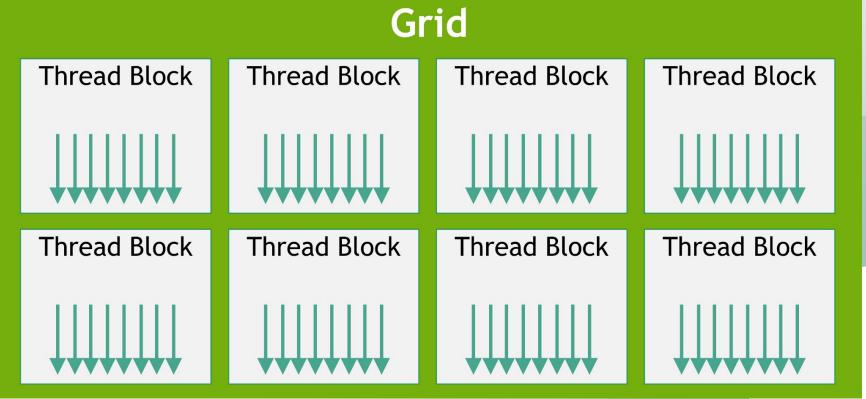
## **Accelerate your Package**



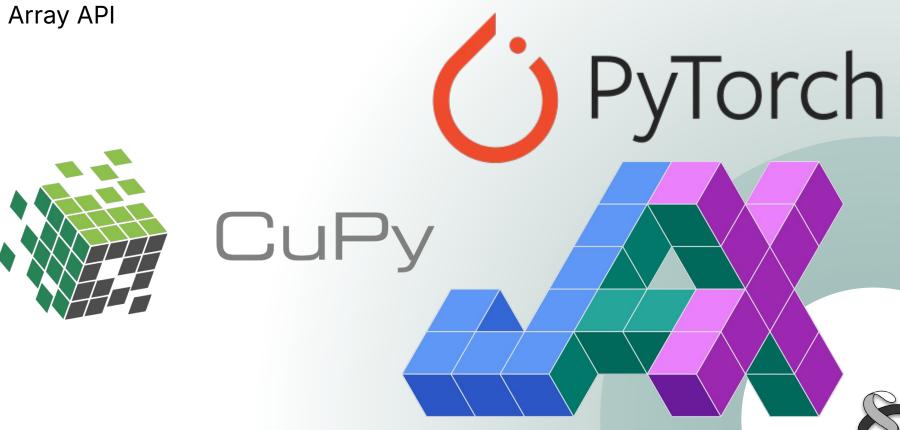




**GPU** 







Rapids



# RAPIDS



### **Get Nerdy**

#### Kernels

- Numba
- CuPy

```
_mul_kernel_csr = r"""
(const int *indptr, {0} *data,
             int nrows, int tsum) {
     int row = blockDim.x * blockldx.x + threadldx.x;
     if(row >= nrows)
        return;
     \{0\} scale = 0.0;
     int start idx = indptr[row];
     int stop idx = indptr[row+1];
     for(int i = start_idx; i < stop_idx; i++)
        scale += data[i];
     if(scale > 0.0) {
        scale = tsum / scale;
        for(int i = start idx; i < stop idx; i++)
          data[i] *= scale;
******
```



#### **Quick Introduction to CSR Matrix**

Sparse Matrix				
10	0	0	0	-2
3	9	0	0	0
0	7	8	7	0
3	0	8	7	5
0	8	0	9	13

