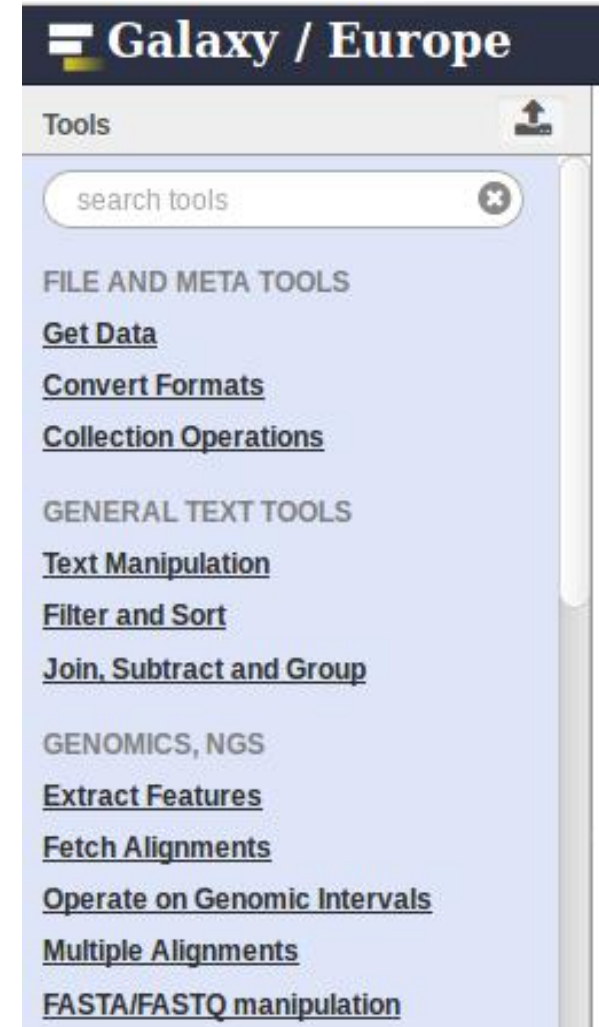# Find similarity in Galaxy tools and predict next tools in workflows
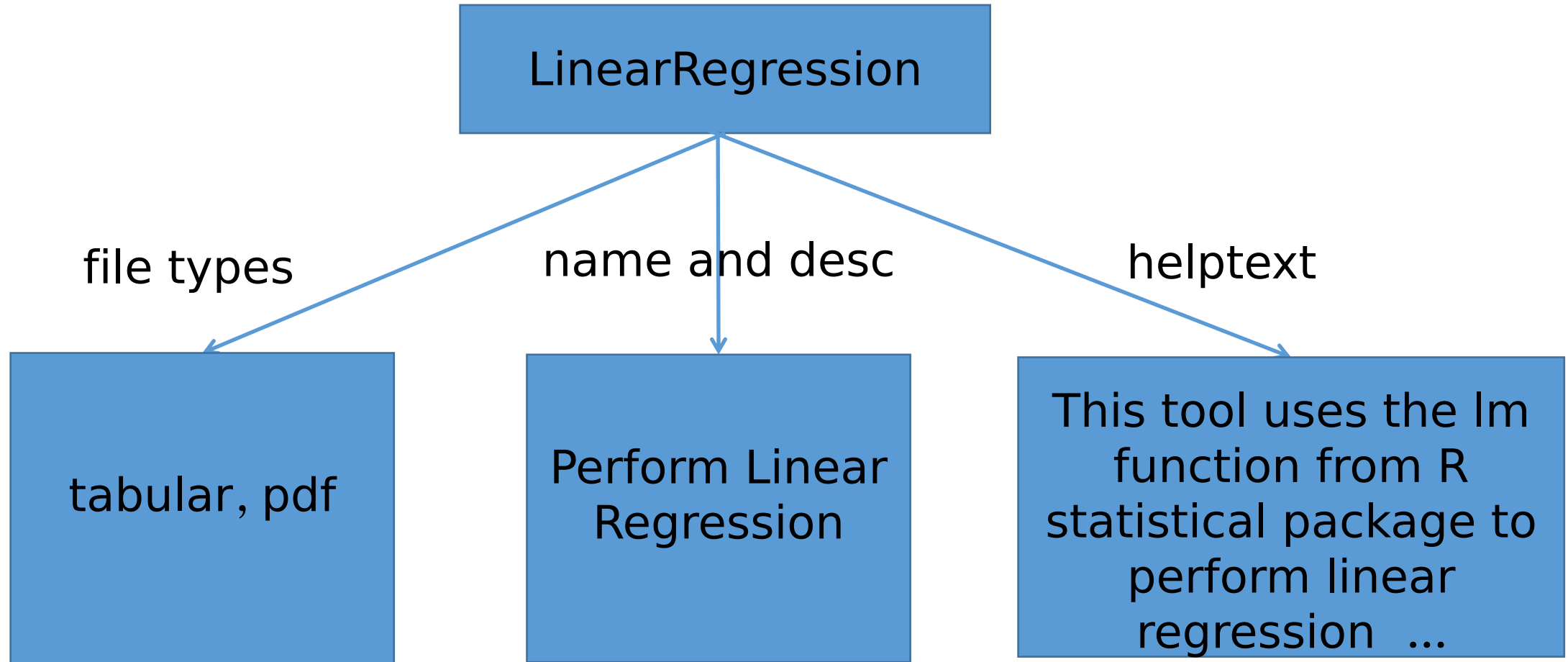
(**Master**'**s thesis**)
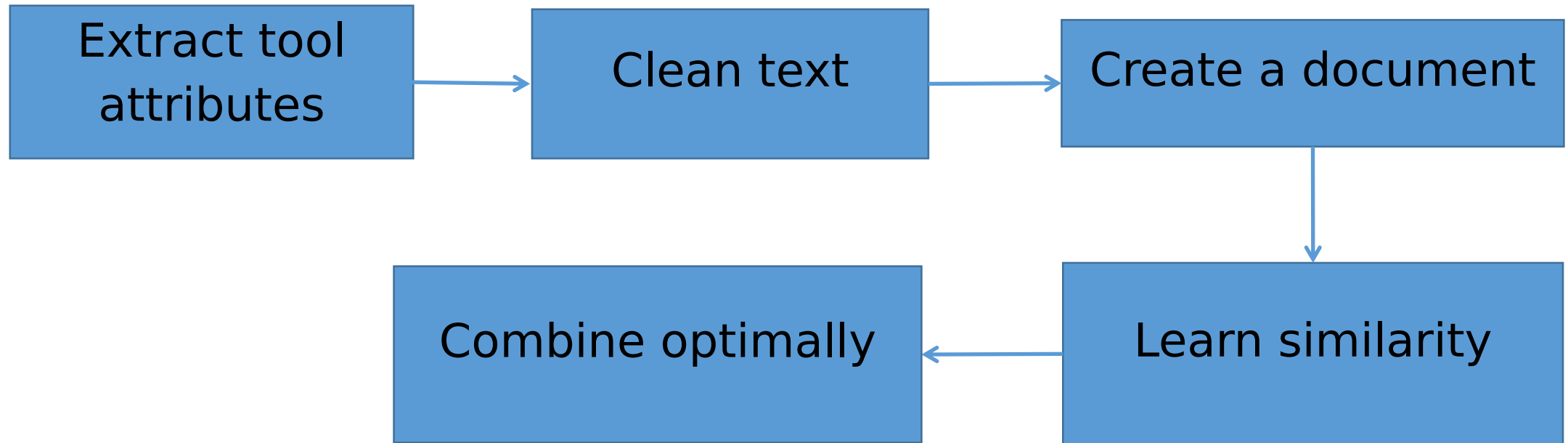
Anup Kumar

# Find similarity in Galaxy tools

Similarity in tools using
machine learning (ML) and natural
language processing (NLP) approaches

# Tool's attributes

# Approach

# Example

| Attributes/ Tools | LinearRegression | LogisticRegression | Similarity |
|---|---|---|---|
| Input, output | 'pdf' , 'tabular' | 'tabular' | ? |
| Name, description | 'regress' , 'linear' , 'perform' | 'logist' , 'regress' , 'perform' | ? |
| Help text | 'regress' , 'assumpt' , 'lm' , 'statist' ,'linear' … | 'vif' ,'regress' , 'glm' ,'car' ,'inflat' , 'function' ,'statist' , 'logist' … | ? |

# Compute similarity

- Jaccard's distance for input/output
- Dense vector for name, description and helptext
- Example:
- ['**regress**', '**linear**', '**perform**'] $= [\ 0.98,\ 0.07,\ ...\ ,\ 0.12\ ]$

# **Similarity matrix** (for name, desc.)

| Tools | LinearRegression | LogisticRegression | BestSubsets Regression | lda_analy |
|---|---|---|---|---|
| LinearRegression | 1 | 0.88 | 0.84 | 0.86 |
| LogisticRegression | .88 | 1 | 0.82 | 0.65 |
| BestSubsetsRegression | 0.84 | 0.82 | 1 | 0.62 |
| lda_analy | ..... | .... | .... | 1 |

# How to combine ?

- $3$ similarity matrices, one for each attribute
- How to combine them ? Take average ?
- Optimal combination, learn weights for each tool
- Similarity:

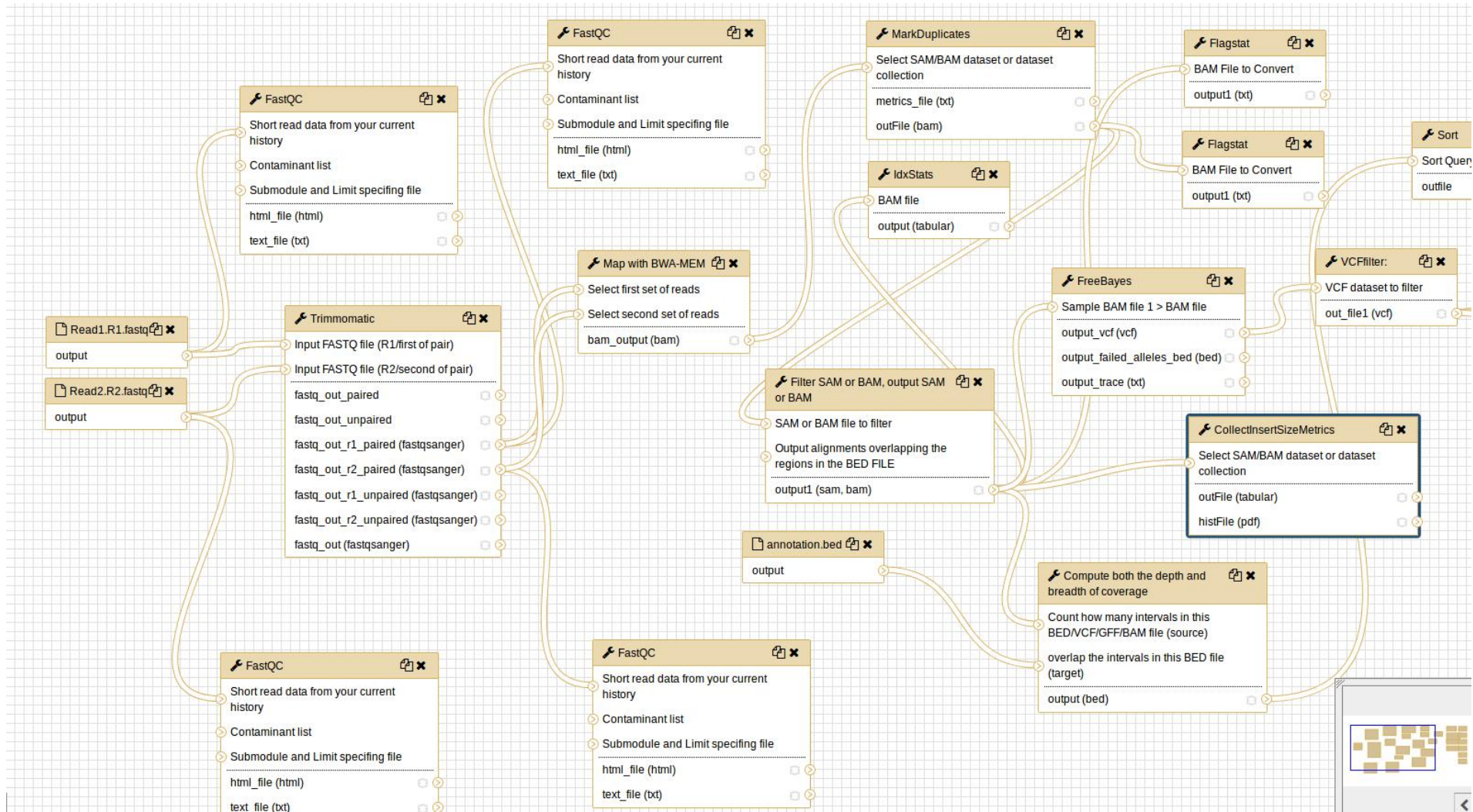$$\arg\max_{(w_i,\ldots,w_n)} \sum_{i=1}^{N} w_i \cdot s_i$$

# Example

- Tool: LinearRegression
- Similarity for input/output: **sim_io** $= [\ 0.33,\ 0.5,\ 1.0,\ .....\ ]$
- Similarity for name, desc: **sim_nd** $= [\ 0.83,\ 0.09,\ 0.005,\ .....\ ]$
- Similarity for helptext: **sim_ht** $= [\ 0.45,\ 0.36,\ 0.001\ .....\ ]$
- Similarity $=$
  $\text{argmax}(w_1\ \times\ \text{sim\_io}\ +\ w_2\ \times\ \text{sim\_nd}\ +\ w_3\ \times\ \text{sim\_ht})$
  $w_1\ +\ w_2\ +\ w_3\ =\ 1$
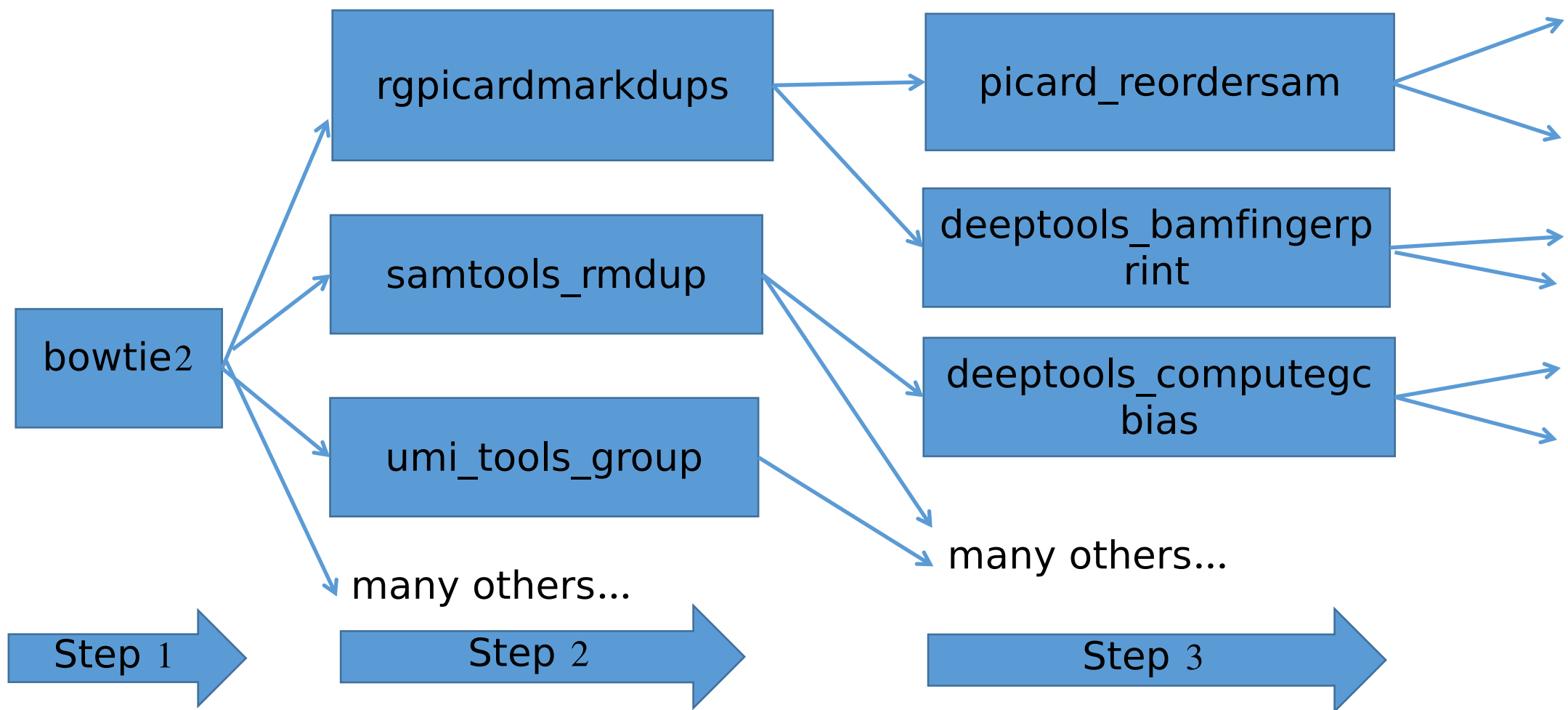
# Visualizer and References

- Static website: results for ~ 1000 tools
- https://rawgit.com/anuprulez/similar_galaxy_tools/master/viz/similarity_viz.html
- https://github.com/anuprulez/similar_galaxy_tools
- https://cs.stanford.edu/%7Equocle/paragraph_vector.pdf
- https://arxiv.org/pdf/1607.05368.pdf

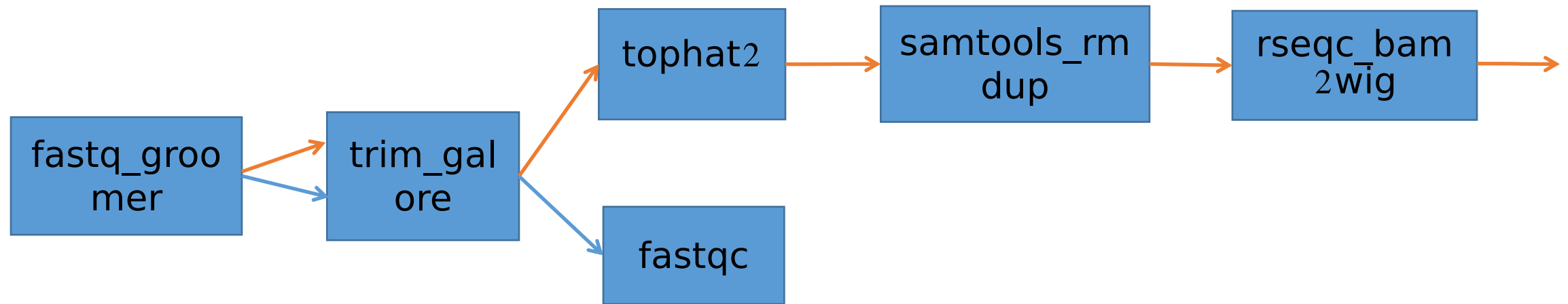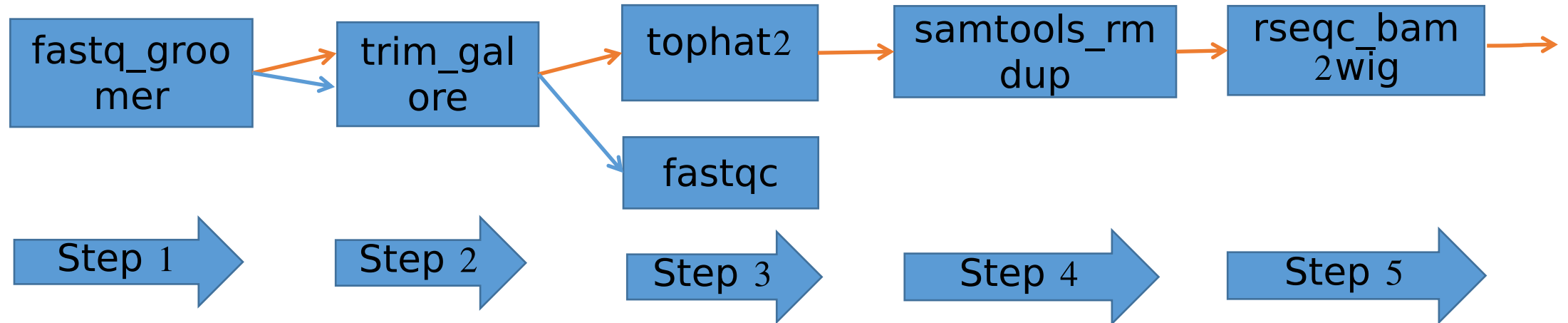# Predict next tools in Galaxy workflows

# Galaxy workflow

# Next tools ?

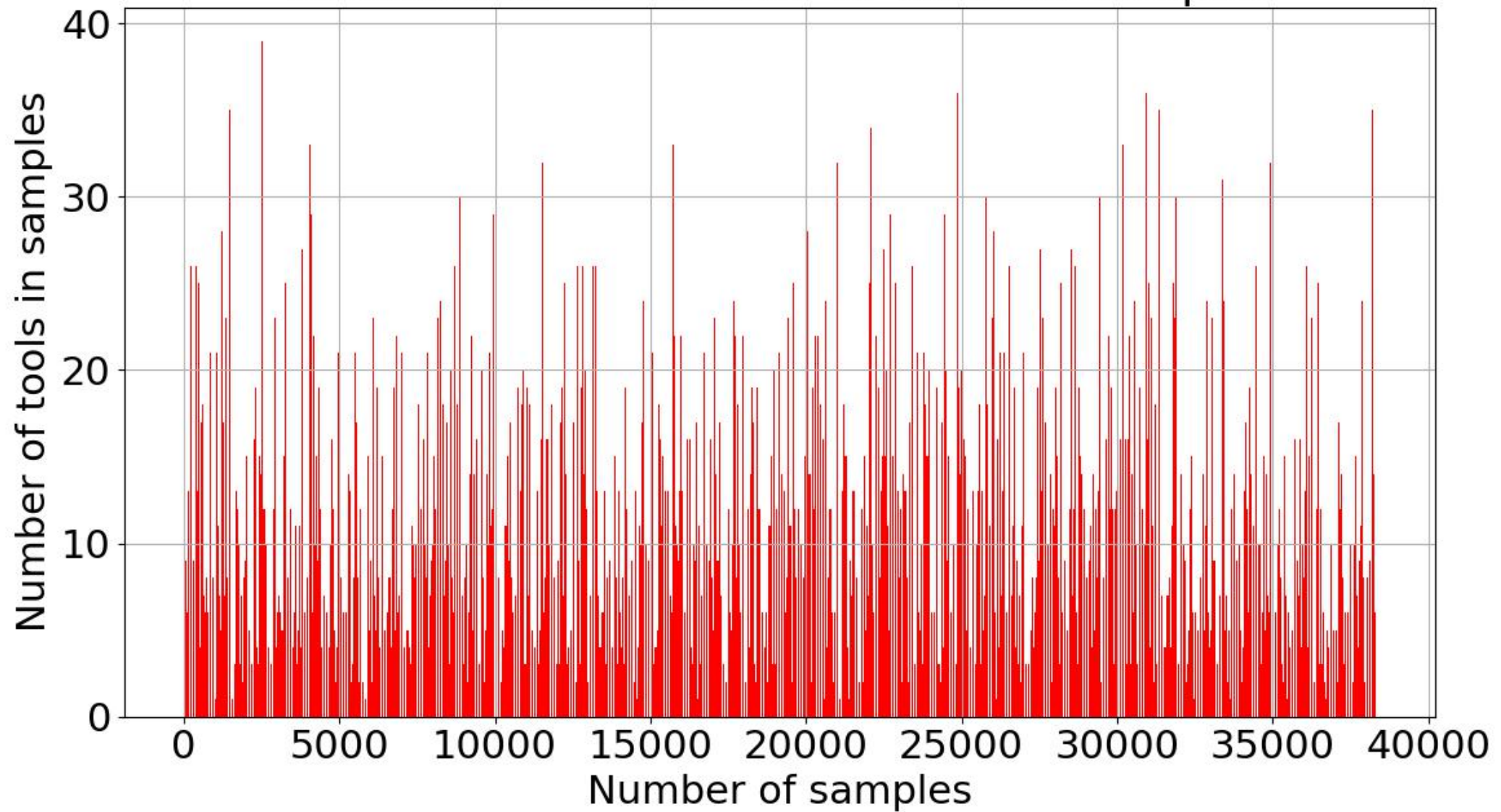# Workflow as a sequence

# Data preprocessing



- fastq_groomer, trim_galore (Step 1)
- fastq_groomer, trim_galore, tophat2, fastqc (Step 2)
- fastq_groomer, trim_galore, tophat2, samtools_rmdup (Step 3)
- fastq_groomer, trim_galore, tophat2, samtools_rmdup, rseqc_bam2wig (Step 4)
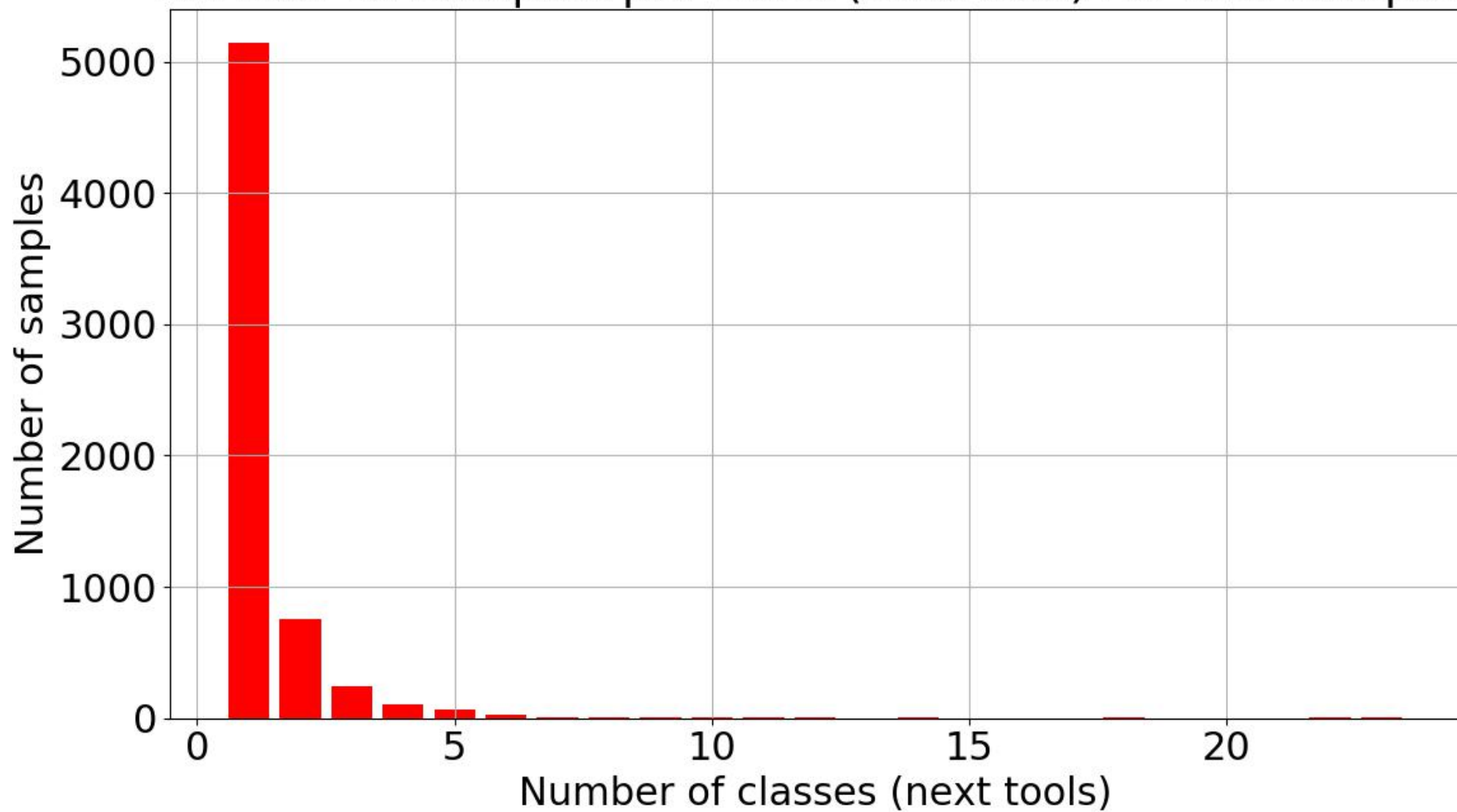
# Data preprocessing

- Assign a unique index to each tool
- {'fastq_groomer': $1$, 'trim_galore': $2$, 'tophat$2$': $3$, 'samtools_rmdup': $4$, 'rseqc_bam$2$wig': $5$, 'fastqc': $6$}
- Training samples:

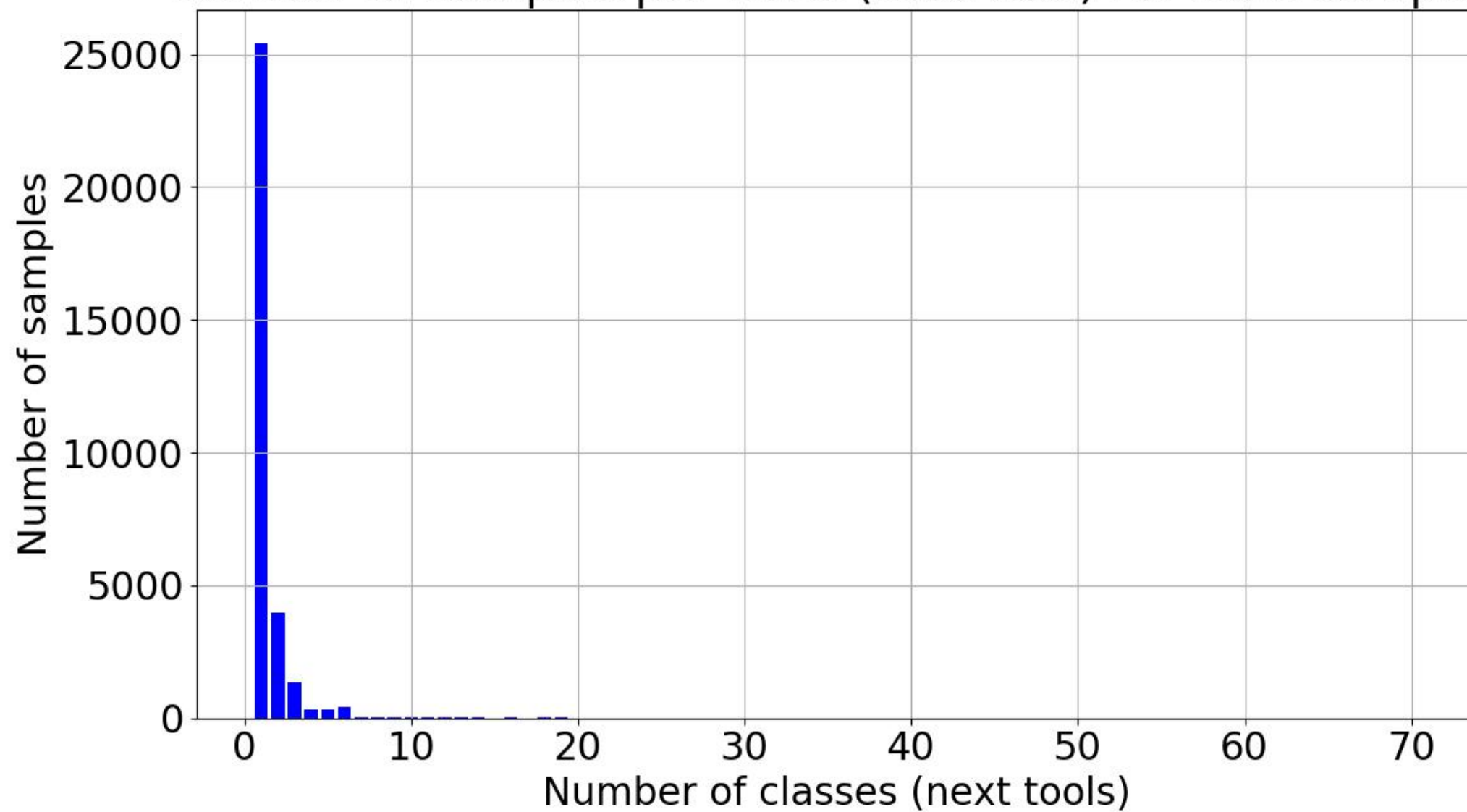| Training sample | Label (next tools) |
|---|---|
| fastq_groomer, trim_galore | tophat$2$, fastqc |
| $1$,$2$ | $3$, $6$ |
| $1$,$2$,$3$ | $4$ |

Distribution of number of tools in samples

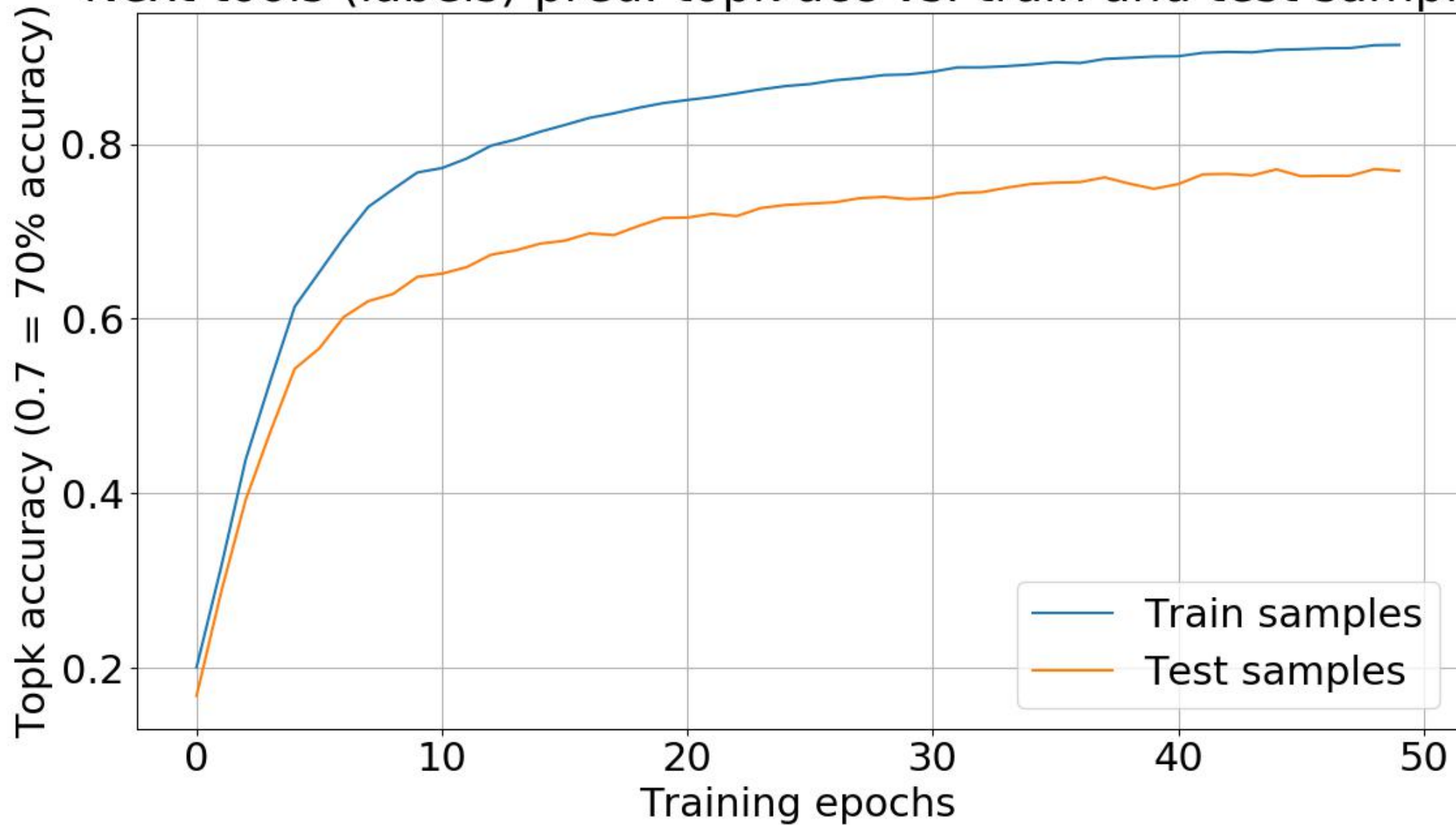Number of samples per class (next tool) for test samples

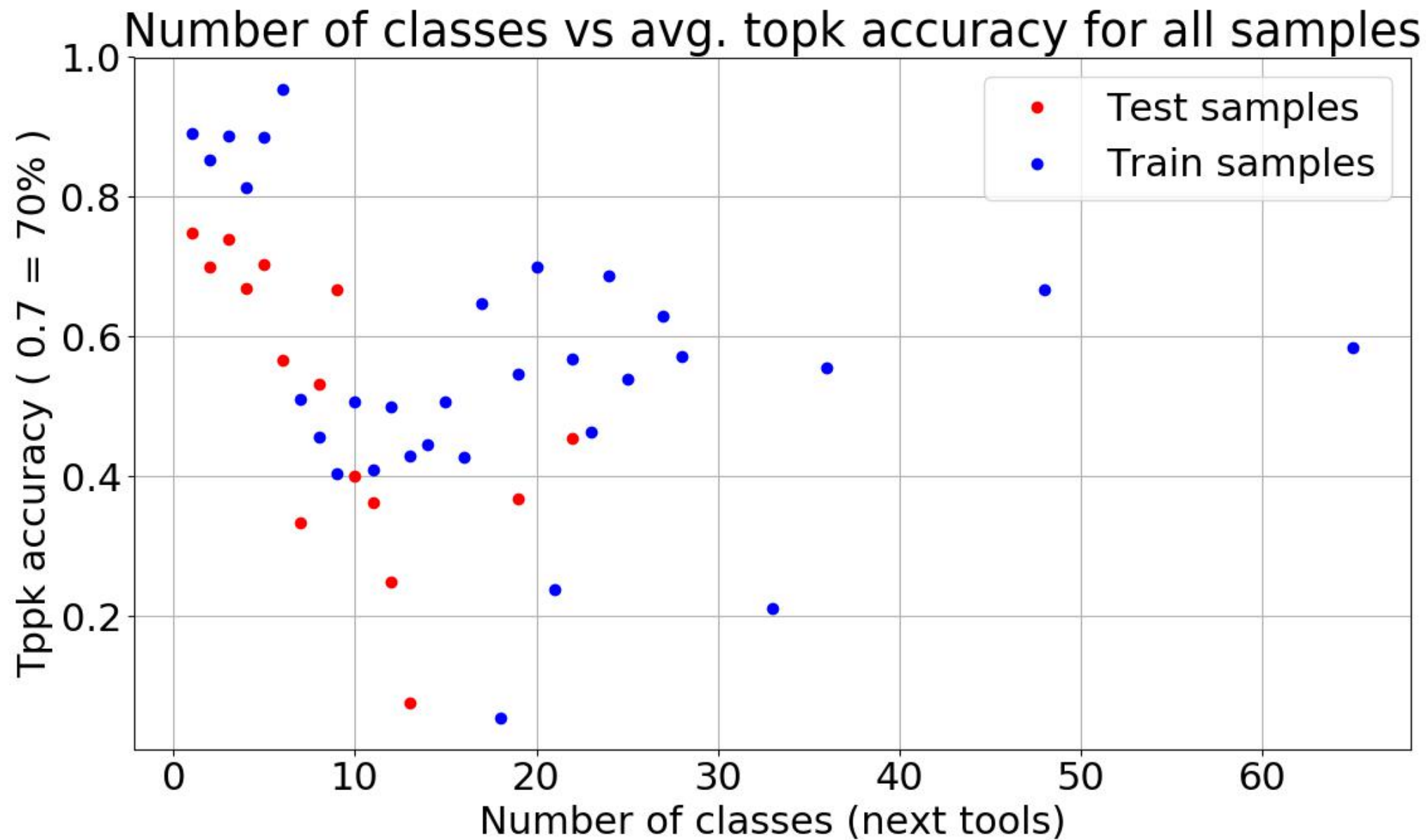Number of samples per class (next tool) for train samples

# Classification

- Multi label, multi class classification
- Long range dependencies samples
- Long short term memory (LSTM) networks
- Topk accuracy
- # of top k tools in actual k next tools ÷ k actual next tools

Next tools (labels) pred. topk acc vs. train and test samples

Number of classes vs avg. topk accuracy for all samples

# Next steps

- Convolution
- Balance the samples
- Different activations
- Compatibility constraint

# References

- https://github.com/anuprulez/similar_galaxy_workflow
- https://arxiv.org/pdf/1511.03677.pdf
- https://arxiv.org/pdf/1604.04573.pdf
- https://arxiv.org/pdf/1506.00019.pdf

# Thank you for your attention