# DEVELOPING TOOLS FOR ANNOTATION AND SINGLE CELL TRANSCRIPTOMICS ANALYSIS

DR WILLIAM HAESE-HILL

WCIP ISAB/RETREAT

31$^{ST}$ MAY 2022

University of Glasgow
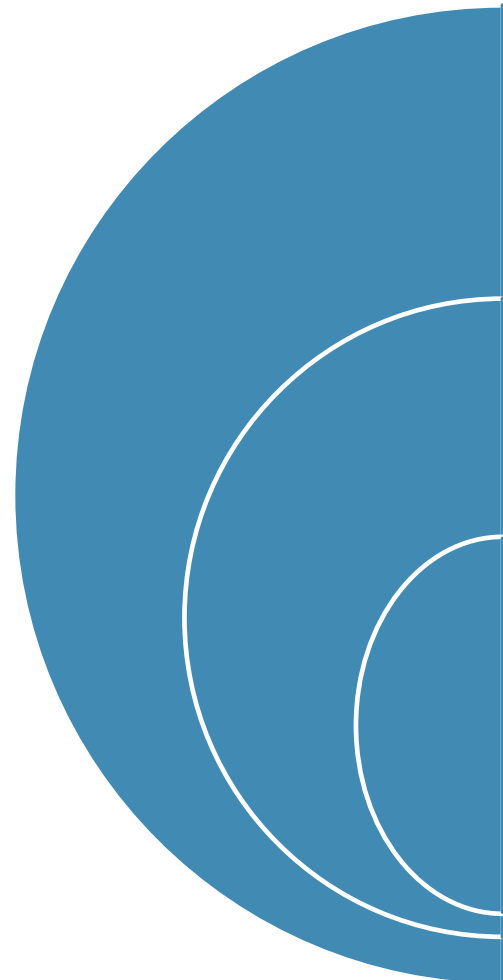
Institute of Infection, Immunity & Inflammation

wellcome centre integrative parasitology

# COMPANION

A GENOME ANNOTATION TOOL FOR MORE THAN JUST PROTISTS

# MOTIVATION

## Why annotate?

- Assign genes with known function to genome assembly
- Recent explosion in assembled sequences
- Aid in drug discovery

## Why automate?

- Numerous sequential processes
- Inputs and outputs same formats
- Improved time / cost efficiency

## Why Companion?

- Only reference-guided annotation tool for eukaryotes
- Scalability
- Potential for larger organisms
- Visualisation outputs
- Already established user base

# PIPELINE

**Pseudochromosome contiguation**
- ABACAS

**Structural annotation**
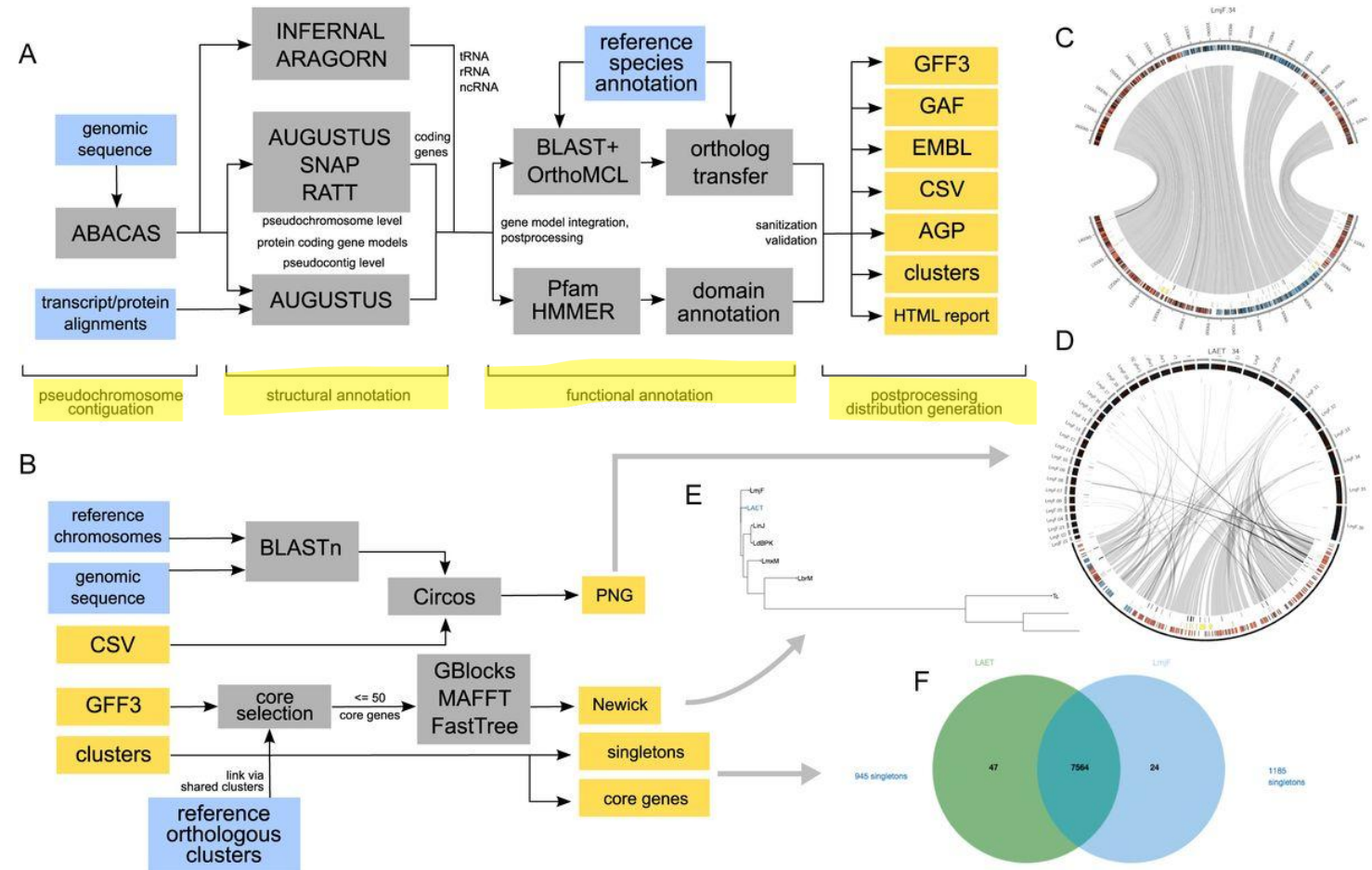- RATT
- SNAP
- AUGUSTUS

**Functional annotation**
- BLAST
- OrthoMCL
- Pfam

**Outputs**

Files: GFF3, EMBL, GAF

Visuals: Orthology, Phylogeny trees, Synteny

# WEB INTERFACE

- Process:
  - Upload a sequence *fasta* file.
  - Select "similar" reference geno̱ from dropdown.
  - Choose various optional proce̱ (e.g. pseudochromosome conti̱ with ABACAS).
  - Submit with (optional) email aḏ for notification.
- Outputs available for up to ̱ months online.
- Download outputs in EBI / GenBank compatible formaṯ.

# NEW FEATURES

CURRENTLY IN DEVELOPMENT

# SCALING FOR LARGER GENOMES

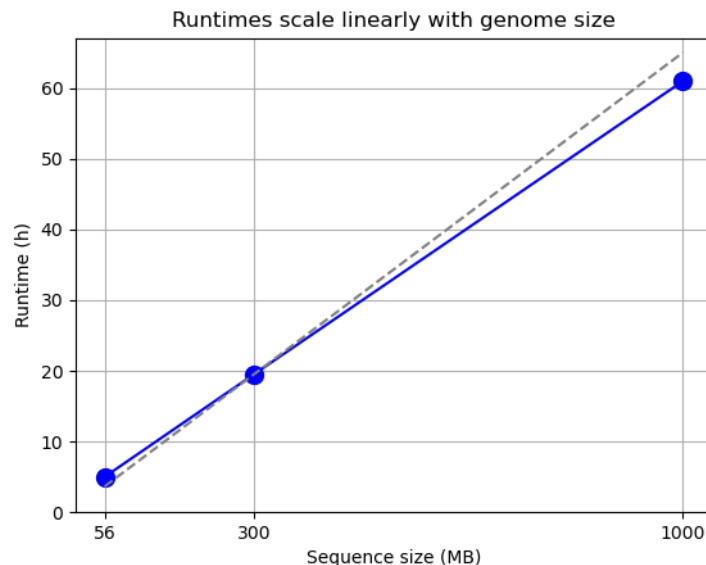References gathered from VEuPathDB projects and pre-compiled.

- In production: protozoa, fungi.
- In development: **VectorBase**
- Exploring: **HostDB**


VectorBase
Bioinformatics Resource for
Invertebrate Vectors of Human Pathogens


HostDB
Pathogen Host Informatics Resources
Release 57
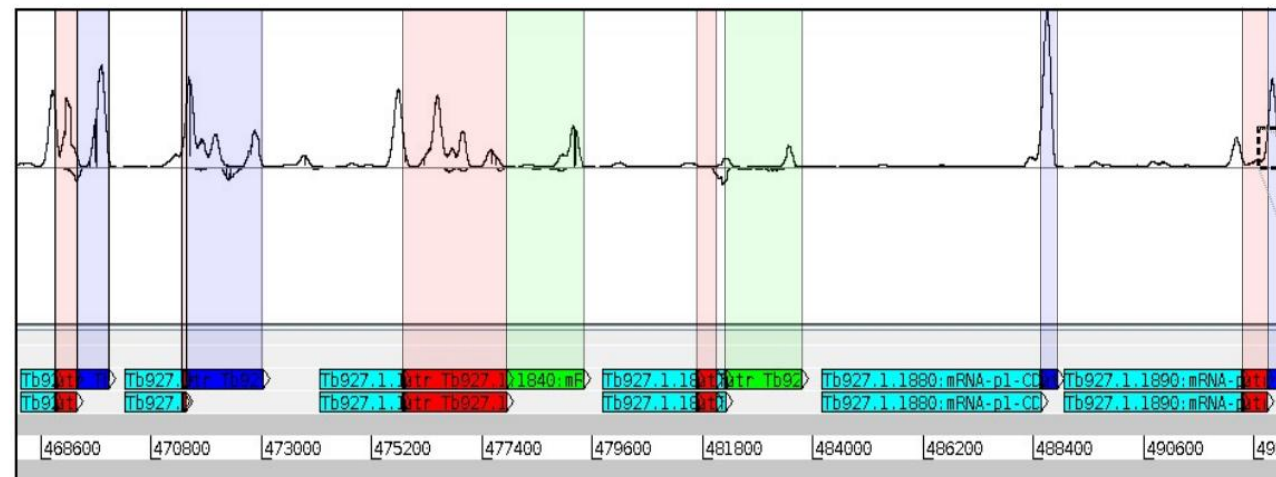21 Apr 2022


Runtimes scale linearly with genome size

- Tested several VectorBase reference genomes
- Runtimes scale linearly
- Issues with RATT for larger genomes
  - Explore fix or alternative tool (such as *Liftoff*)

Alaina Shumate, Steven L Salzberg, **Liftoff: accurate mapping of gene annotations**, *Bioinformatics*, Volume 37, Issue 12, 15 June 2021, Pages 1639–1643, https://doi.org/10.1093/bioinformatics/btaa1016

# UTR ANNOTATION

- **peaks2utr**: stand-alone 3' UTR annotation Python tool

- How to incorporate:

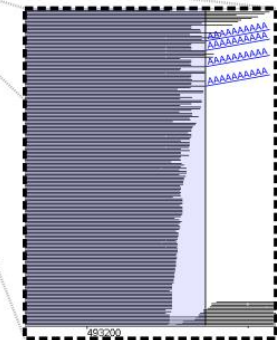  - Apply to select reference annotations - or to output annotations?



*T. brucei TREU927*

# FUTURE ADDITIONS

- Improved efficiency / job concurrency.

- Expanded reference set.

- Dynamic reference updates.

- Improved visualisations (e.g. Apollo).

- Full project submission for EBI / GenBank

- Additional functionality to increase richness of annotations.

# CELL ATLAS – cellxgene
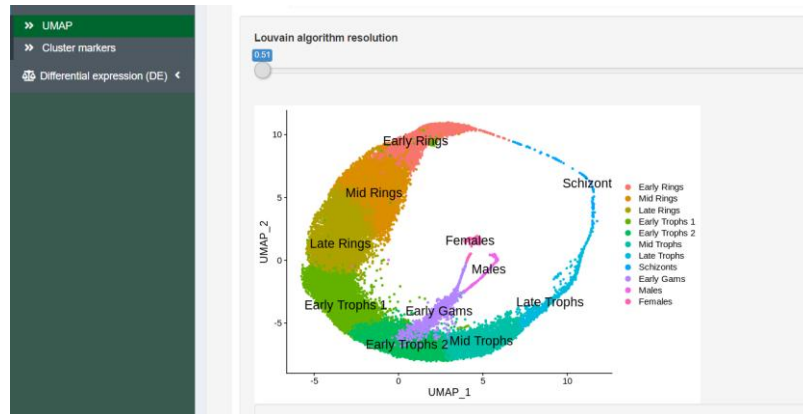
RESPONSIVE WEB INTERFACE FOR VISUALISING SINGLE CELL DATA

# WHAT IS CELLXGENE?

- Interactive single cell (SC) analysis framework.

- Developed by Chan Zuckerberg Initiative for the Human Cell Atlas.

- Predominantly Python with Flask app interface.

- Simple to install and launch from command line.

- Third-party tools for hosting (cellxgene-gateway) and enhancing featureset (cellxgene_VIP, excellxgene).
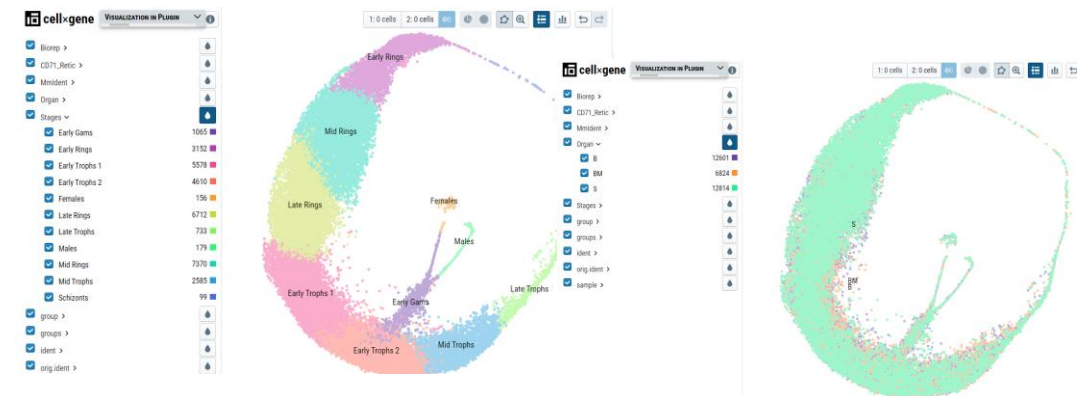
# COMPARISON

## UofG III - Cell Atlas

- R/Shiny interface.

- Must be maintained by in-house developers at UofG.

- Each dataset requires a bespoke approach.
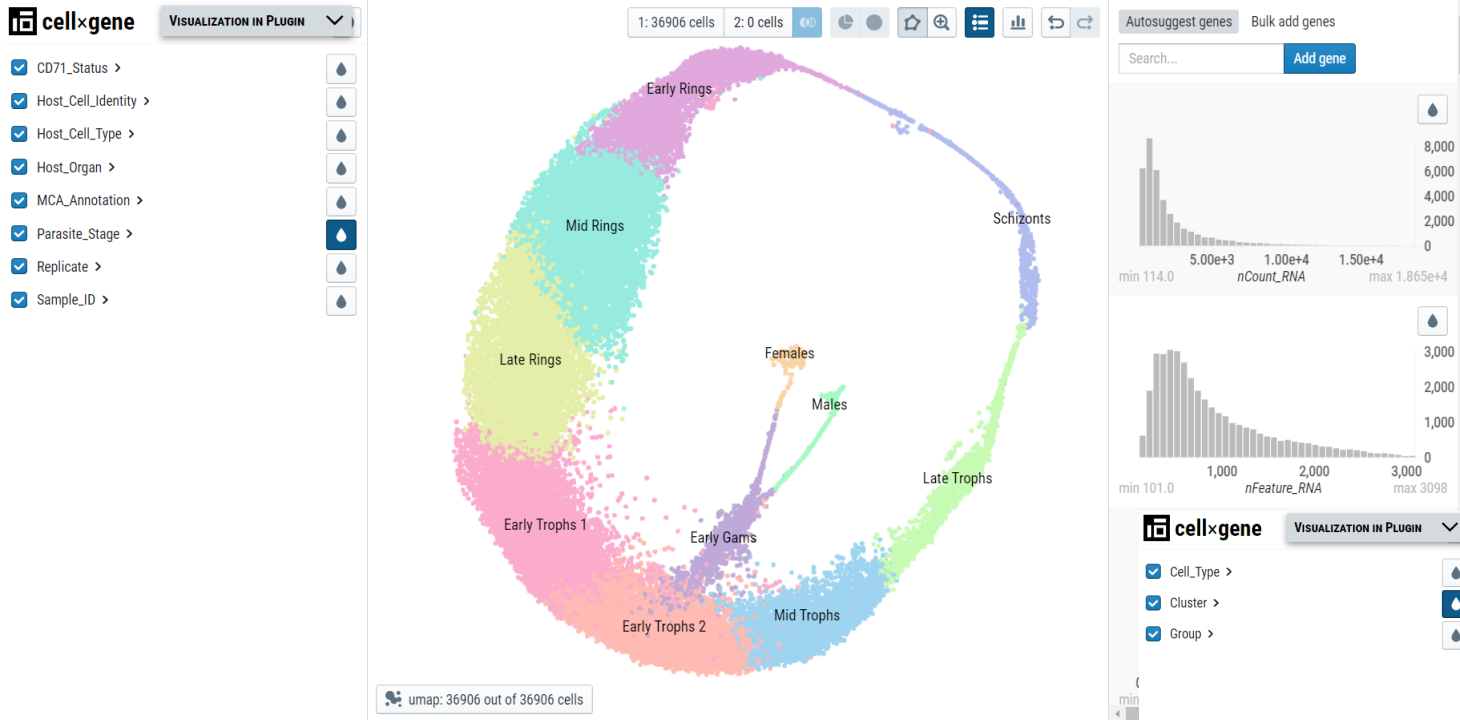
- Manually defined plotting groups.



- Unresponsive interface effects user experience.

## Cell Atlas - cellxgene

+ Python/Flask interface.

+ Regularly maintained GitHub repos.

+ All functionality included out-of-the-box for each dataset.
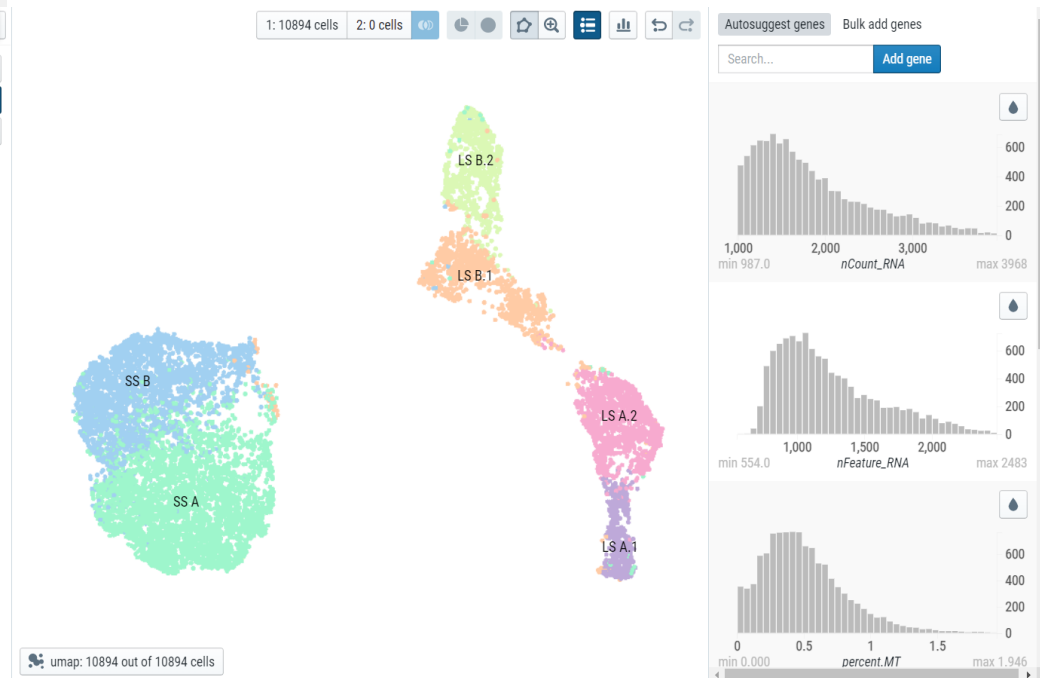
+ Plotting groups dynamic and inferred from dataset.



+ Once dataset loaded, interface is snappy.

+ Includes all the same features, and then some.

https://cellatlas-cxg.mvls.gla.ac.uk/Tbrucei/

*Trypanosoma brucei*

*Plasmodium berghei*

https://cellatlas-cxg.mvls.gla.ac.uk/Pb/

# ADDITIONAL FEATURES

CURRENTLY IN DEVELOPMENT

# VEUPATHDB INTEGRATION


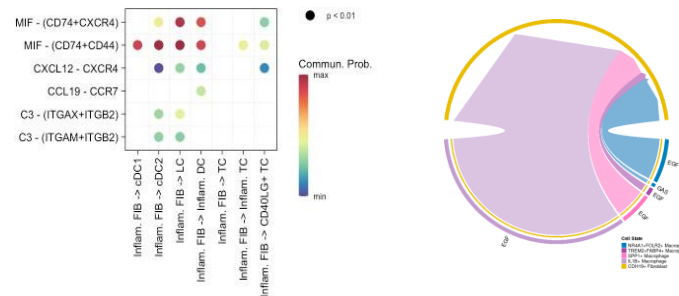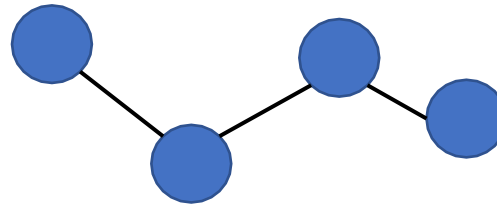
- Gene detection algorithm
  – link to various VEuPathDB sub-sites

# EXPLORATION AND VISUALIZATION OF CELL-CELL INTERACTIONS*

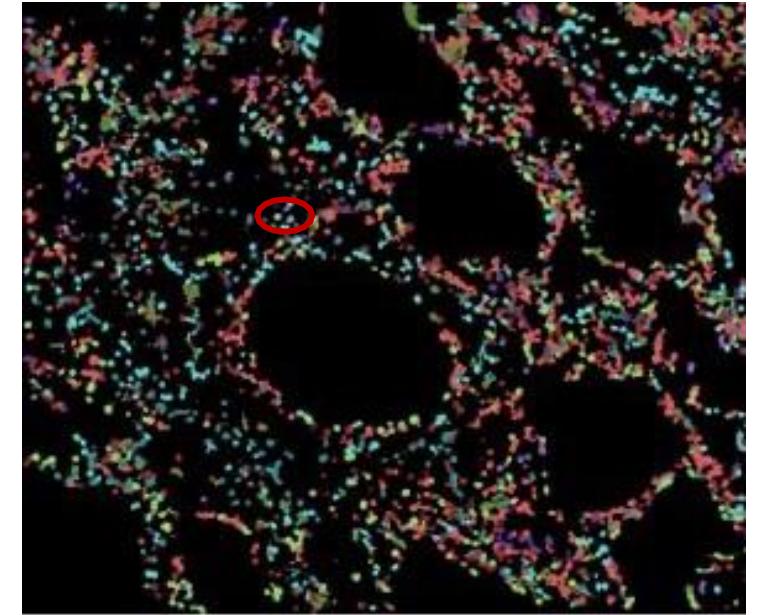Combined projection of interacting cells

Single cell-RNA sequencing data

User customised plots

Spatial data

* PhD project of Olympia Hardy

'Cell Type A'

'Cell Type B'

RESEARCH INTO INFLAMMATORY ARTHRITIS CENTRE
VERSUS ARTHRITIS

# ACKNOWLEDGEMENTS

Glasgow Bioinformatics Summerschool

Application Deadline: 07/07/2022



| | | | |
|---|---|---|---|
| Thomas Otto | Ross Laidlaw | Lucy MacDonald | Emma Briggs |
| Kathryn Crouch | Fiona  Achcar | Eva Crespo | Domenico Somma |
| Edward Agboraw | Alex Pancheva | Lauren Galloway | Theodore Simakou |
| Olympia Hardy | John Cole | Collins Morang'a | Katie Chapple |
| Scott  Arkison | | | |

University of Glasgow

wellcome centre integrative parasitology