# Bioinformatics Application Note:

# CONTEXT – A Phylogenomics Dataset Browser

Joe Parker[1,2] and Stephen J. Rossiter[2]

1. Kitson Consulting, Bristol, UK; Present address: Jodrell Laboratory, Royal Botanic Gardens, Kew, UK

2. School of Biological and Chemical Sciences, Queen Mary University of London, UK.

**Word count:** <u>1,035</u>

Corresponding Author:
Joe Parker
Jodrell Laboratory,
Royal Botanic Gardens, Kew,
TW9 3DS, UK
Tel. +44 20-8332-5063
Fax +44 20-8332-5197
joe.parker@kew.org

Project email: joe+CONTEXT@kitson-consulting.co.uk

## *Abstract*

**Summary.** Quality control (QC) in large phylogenomic datasets is a key requirement for reliable and reproducible research in evolution, adaptation, speciation and taxonomy. CONTEXT is a browser for high-throughput visualisation and comparative QC of phylogenomic datasets, consisting of a Java API and an executable binary jarfile with graphical user interface (GUI). The tool allows users to rapidly and easily visualise thousands of multiple sequence alignments and hundreds of phylogenies using a GUI to identify outliers which could affect downstream analyses. CONTEXT calculates a variety of downstream statistics on alignments and phylogenies including entropy, informativeness, imbalance, signal:noise and size.

**Motivation.** Comparative genomics studies have become increasingly common, but these analyses are sensitive to the quality and heterogeneity of input datasets (multiple sequence analyses and phylogenies). Currently few tools exist to readily compute descriptive statistics, or to visualise large numbers of input datasets. CONTEXT is a phylogenomics dataset browser which facilitates these analyses in a lightweight application. It allows any user to rapidly visualise, inspect, score, and sort input datasets to identify outlying datasets which may need additional processing, filtering, or masking from further analyses.

**Results.** The application has been successfully implemented on a variety of infrastructures. A variety of common input data formats including FASTA, Phylip/PAML, Nexus, and Newick conventions are automatically read and parsed.

**Availability and implementation.** The API is implemented in native Java code, available online at https://github.com/lonelyjoeparker/qmul-genome-convergence-pipeline. The executable binary can be downloaded at https://github.com/lonelyjoeparker/qmul-genome-convergence-pipeline/tree/master/trunk/bin. The project page is at https://github.com/lonelyjoeparker/qmul-genome-convergence-pipeline/blob/master/CONTEXT.md

**Contact.** joe.parker@kew.org

## *Introduction*

In the last decade, 'phylogenomic' analyses of sequence data (typically, phylogenetic analyses comprising genome-scale data, e.g. thousands of coding sequences) have become increasingly common. They have provided a powerful means to highlight evolutionary insights in animals (Tsagkogeorga *et al.*, 2013; McCormack *et al.*, 2013), plants (Zhao *et al.*, 2016), fungi (Dentinger *et al.*, 2016) and prokaryotes (Jungblut *et al.,* 2016). As well as pure phylogenetic inference, the greatly expanded data these analyses drive permit detection of subtle but key genomic effects such as adaptive molecular convergence (Parker *et al.*, 2013) and incomplete lineage sorting (Scornavacca & Galtier, 2016), among others. The trend towards using genomic big data to answer evolutionary questions will surely only accelerate.

However, the multiplicity of sequencing methods, annotation and orthology protocols, multiple sequence aligners, and phylogeny inference tools available mean that comparing the inputs to downstream analyses (typically, multiple sequence alignments and phylogenies) is essential. In particular, between-workflow differences can lead to widely differing results, even for identical data (Blackburne & Whelan, 2013). For this reason, robust quality control procedures should be applied to dataset curation and phylogeny inference – with repeatability an essential criterion. Unfortunately, many labs' workflows have yet to meet this demand; typically only limited comparative QC analysis is performed on sets of alignments and phylogenies. In part, this is because, although certain descriptive statistics are available, they are not routinely calculated by upstream (genomics) packages such as assemblers or annotation tools. In many cases, inspection by eye remains the primary means of error detction – and even then, often only used *post hoc* to troubleshoot outliers in downstream analyses. This situation is clearly untenable in the phylogenomic era.

CONTEXT (for 'COmparative Nucleotides and Trees EXploration Tool') is designed to improve this situation. To improve the objectivity by which alignments and phylogenies are selected for further analyses, various helpful statistics are rapidly calculated to aid dataset QC. Their summaries are also calculated at the dataset level, and simple visual plots can be used to identify outliers. Meanwhile, the graphical user interface provides a fast method by which thousands of alignments and phylogenies can be loaded and inspected by eye if the user wishes.

## *Features and implementation*

The API elements contain resources for phylogenomics such as input/output and parsing utilities; trimming, pruning and validation methods for alignments and phylogenies; statistics for evaluating alignments, phylogenies, likelihood fits and dN/dS values; UI elements including two main GUI platforms; post-processing including linear regression and descriptive parametric statistics on large distributions of small floating-point numbers. Implemented alignment statistics include (for nucleotide and amino-acid sequences): number of taxa; length (excluding and including gaps and invariant sites); longest sequence without gaps; mean sitewise entropy (after Shannon (1948)) and longest run of non-zero entropy. Phylogeny statistics (after Felsenstein, 2004) include number of tips; tree and root length and height; imbalance (Colless, 1982); cherry-count (Steel & McKenzie, 2001); signal:noise (aka internal branch lengths:total tree length Phillips & Penny, 2001) and external:internal branch length ratio.

## *Evaluation*

Example usage statistics shown in Table 1.

| | OS | Arch | CPU type, clock GHz | cores | RAM Gb | HDD Gb |
|---|---|---|---|---|---|---|
| pandanus | Ubuntu 16.04 LTS / Biolinux | i686 | Xeon E5620 @ 2.4 | 4 | 33 | 1000 @ ATA 7200rpm |

| | | | | | | |
|---|---|---|---|---|---|---|
| 2 | Ubuntu 16.04 LTS / MATE | ARM | ARMv7 @ 0.9 | 1 | 1 | 8 @ SD |
| toshiba | Windows 7 | x64 | Core i7 @ 2.4 | 4 | 3 | 128 @ SSD |
| Tower | Ubuntu 16.04.02 LTS | X64 | 8x Intel(R) Xeon(R) CPU E3-1245 v5 @ 3.50GHz | 8 | 65 | 1.0Tb SSD, 4.0Tb HDD |
| MBP | Mac OSX 10.9.5 | x64 | Core i7 @ 2.2 | 4 | 8 (1333MHz DDR3) | 250 @ SSD |
| EC2 m4.10xlarge | Ubuntu 16.04 LTS | x64 | Xeon E5-2670 @ 2.5 | 16 | 122 | 320 @ SSD |
| EC2 | Ubuntu !6.04 LTS | x64 | Xeon E5-2680 @ 2.8 | 32 | 60 | 2x320 @ SSD |

**Table 1**

## *Roadmap and versioning*

CONTEXT is currently supplied at **Version 0.8.4 prerelease.**

## *Acknowledgements*

## *Figures / data / tables*

**Table 1: Example system resource usage.** The RAM usage (in megabytes) and average load time (initialisation to completion) of the CONTEXT under a variety of test computer architectures and input datasets. Benchmarking scripts and analyses are available online at https://github.com/lonelyjoeparker/qmul-genome-convergence-pipeline/benchmarking

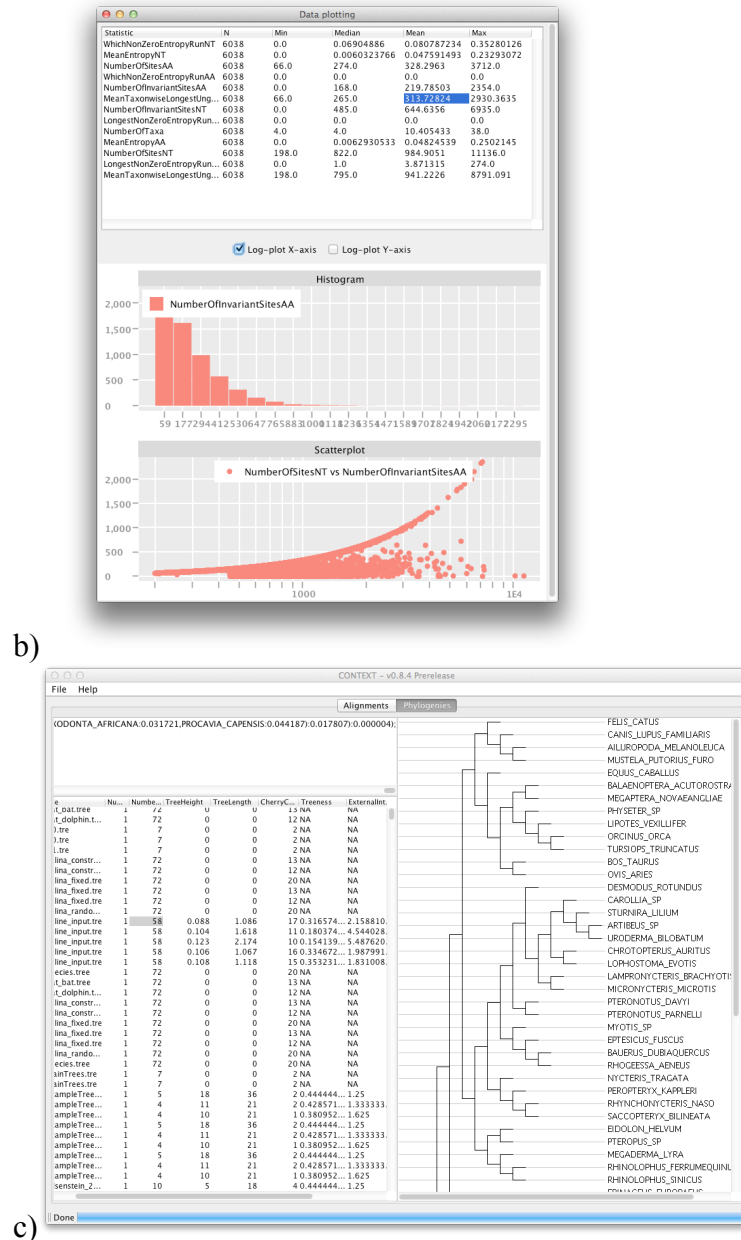| Test case | Mac OSX 10.9.5, 2.2GHz core i7, 8Gb 1333MHz DDR3 RAM, 250 Gb SSD. | Tower Ubuntu 16.04.02, 8x Xeon E3-1245v5 3.5GHz, 65Gb RAM, 1.0Tb SSD/4.0Tb HDD | Windows 7, 2.0GHz core i7 CPU, 8Gb RAM, 500Gb SSD |
|---|---|---|---|
| 692 Nucleotide alignments, 7 taxa, 78-2766nt | *6.8* | *5.0* | |
| 2,326 Nucleotide alignments, 14-22 taxa, 450-11136nt | *21.7* | *19.0* | |
| 3,656 Nucleotide alignments, 10 taxa, 99-6366nt | *25.1* | *23.0* | |

**Table 1**

**Table 2: Load time benchmarking.** The RAM usage (in megabytes) and median load time (initialisation to completion) of the CONTEXT for scaled task sizes, on two representative systems; a multipurpose Linux bioinformatics workstation in general use mode (other applications active simultaneously) and an Apple laptop running in dedicated mode (no other user tasks active). Values shown wall-clock time in seconds to initialise CONTEXT, load the specified dataset and calculate statistics, and exit (median values of five measurements).

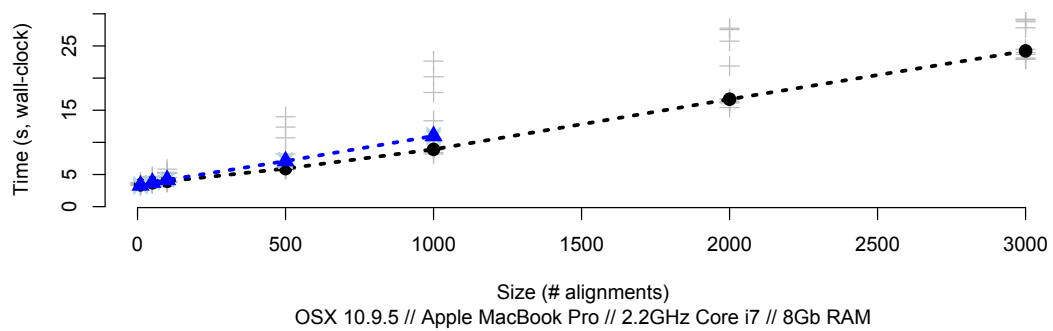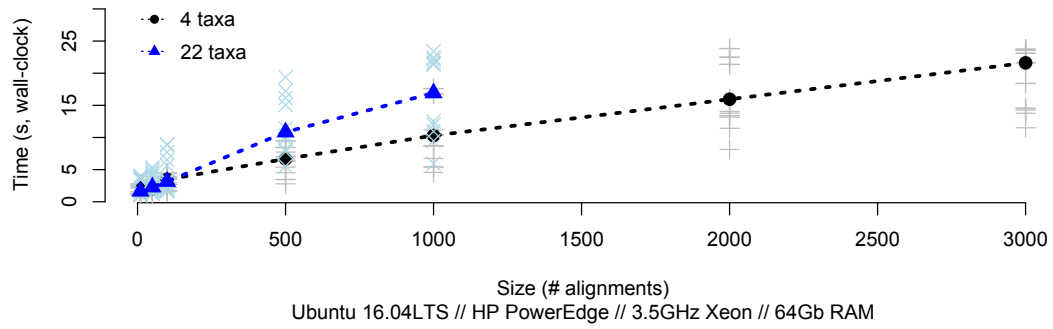| System | Tower Ubuntu 16.04.02, 8x Xeon E3-1245v5 3.5GHz, 65Gb RAM, 1.0Tb SSD/4.0Tb HDD | | | Mac OSX 10.9.5, 2.2GHz core i7, 8Gb 1333I DDR3 RAM, 250 Gb SSD. | | |
|---|---|---|---|---|---|---|
| **Taxa à** **Alignments** ∨ | **4** | **22** | **200** | **4** | **22** | **20** |
| **10** | 2.18 | 1.62 | | 3.39 | 3.30 | |
| **50** | 2.05 | 2.28 | | 3.71 | 3.74 | |
| **100** | 3.46 | 3.14 | | 3.94 | 4.16 | |
| **500** | 6.67 | 10.84 | | 5.92 | 7.09 | |
| **1000** | 10.30 | 16.92 | | 8.92 | 10.99 | |
| **2000** | 15.95 | | | 16.72 | | |
| **3000** | 21.61 | | | 24.23 | | |

**Table 2**

**Figure 1: CONTEXT schematic**. The schematic logic flow of the phylogenomic dataset browser is shown with descriptions of key analysis steps, in flow diagram format.
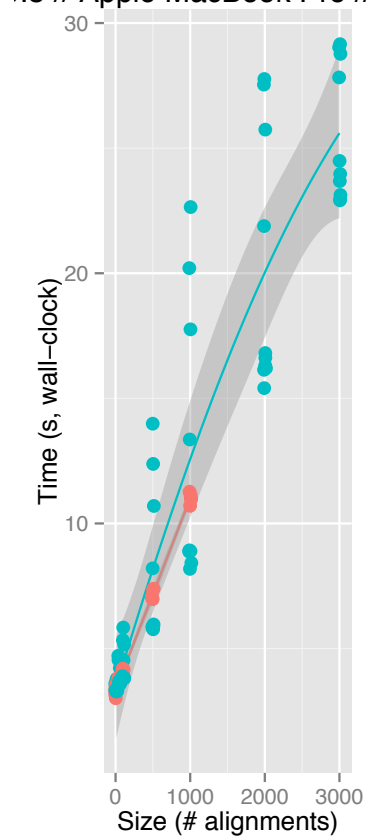
**Figure 2: Screenshots showing visualisation of example datasets:** (a) The alignment browser, showing 692 multiple sequence alignments together with statistics; (b) The phylogeny browser, showing phylogenies and statistics; (c) The plotting display, showing summary statistics, histogram and bivariate data (logged, in this case)
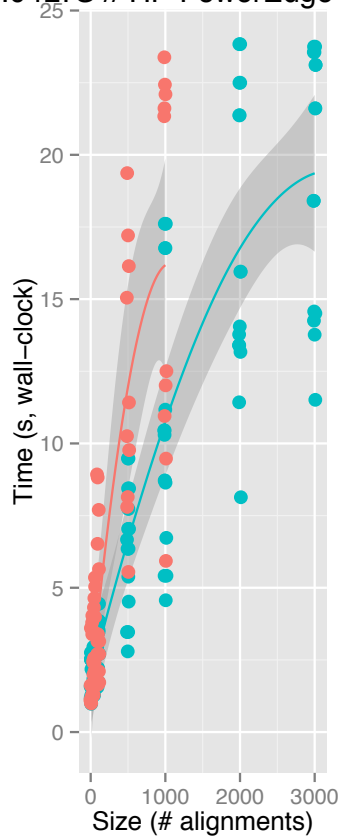


a)

b)



c)

**Figure 3: Performance of CONTEXT in preload task, 50-3000 alignments:** (a) Load times of the Linux workstation in mixed-use mode. 4-taxa alignment tasks shown in black dots, 22-taxa load times as blue triangles. Measurements shown in light points with medians (*n*=5) shown as solid-point; (b) Load times of the Apple laptop in dedicated mode. Labels as above.

# Performance



Size (# alignments)
Ubuntu 16.04LTS // HP PowerEdge // 3.5GHz Xeon // 64Gb RAM



Size (# alignments)
OSX 10.9.5 // Apple MacBook Pro // 2.2GHz Core i7 // 8Gb RAM

# References

Blackburne, B.P. & Whelan, S. (2013) Class of multiple sequence alignment algorithm affects genomic analysis. *Mol. Biol. Evol.* 30(**3**):642-53.

Dentinger, B.T.M., Gaya, E., O'Brien, H., Suz., L.M., Lachlan, R., Díaz-Valderrama, J.R., Koch, R.A. & Aime, C. (2016) Tales from the crypt: genome mining from fungarium specimens improves resolution of the mushroom tree of life. *Bot. J. Linn. Soc* **117**(1):11–32. DOI: 10.1111/bij.12553

Zhao, L., Li, X. , Zhang, N., Zhang S-D., Yi, T-S., Ma, H., Guo, Z-H. & Li, D-Z. (2016) Phylogenomic analyses of large-scale nuclear genes provide new insights into the evolutionary relationships within the rosids. *Mol. Phylogenet. Evol.* **105**:166–176.

Parker, J.*, Tsagkogeorga, G.*, Cotton, J.A., Liu, Y., Provero, P., Stupka, E. & Rossiter, S.J. (2013) Genome-wide signatures of convergent evolution in echolocating mammals. *Nature* **502**(7470)**:**2280231. doi:10.1038/nature12511. *These authors contributed equally to this article.

Tsagkogeorga, G.*, Parker, J.*, Stupka, E., Cotton, J.A., & Rossiter, S.J. (2013) Phylogenomic analyses elucidate evolutionary relationships of the bats (Chiroptera) *Curr. Biol.* **23**(22):2262-2267.

Scornavacca, C. & Galtier, N. (2016) Incomplete Lineage Sorting in Mammalian Phylogenomics. *Syst. Biol.* Xxxxx.

McCormack JE, Harvey MG, Faircloth BC, Crawford NG, Glenn TC, Brumfield RT. (2013) A phylogeny of birds based on over 1,500 loci collected by target enrichment and high-throughput sequencing. *PLoS One.* **8**(1):e54848. doi: 10.1371/journal.pone.0054848.

Anne D Jungblut, Ian Hawes, Tyler J Mackey, Megan Krusor, Peter T Doran, Dawn Y Sumner, Jonathan A Eisen, Colin Hillman, Alexander K Goroncy (2016) Microbial mat communities along an oxygen gradient in a perennially ice-covered Antarctic lake *Appl. Env. Microbiol.* **82**(2):620-630

Shannon (1948)

Felsenstein (2004)

Colless, (1982)

Steel & McKenzie (2001)

Phillips & Penny (2001)