# NGSANE
## A HPC Processing Framework for Terabyte-scale Sequencing Data

GARVAN INSTITUTE

Fabian A. Buske[1], Susan J. Clark[1], Denis C. Bauer[2]

[1] Cancer Epigenetics Program, Cancer Research Division, Garvan Institute of Medical Research, Kinghorn Cancer Centre, Sydney, NSW. 2010 Australia
[2] Division of Computational Informatics, CSIRO, Sydney, NSW. 2113 Australia

CSIRO

The initial steps in the analysis of Next Generation Sequencing (NGS) data can be automated by way of software 'pipelines'. However, individual components deprecate rapidly due to evolving technology and analysis methods, often rendering entire versions of production informatics pipelines obsolete. Constructing pipelines from Linux bash commands enables the use of hot swappable, modular components as opposed to the more rigid program-call wrapping by higher level languages, as implemented in comparable published pipelining systems. Here we present NGSANE, a Linux-based, HPC-enabled framework that minimises overhead for set-up and processing of new projects yet maintains full flexibility of custom scripting when processing raw sequence data.

### HPC & parallel execution

NGSANE supports Sun Grid Engine (SGE) and Portable Batch System (PBS) job scheduling and can be operated in different modes for development and production thus enabling efficient and flexible processing of NGS data. HPC job partitioning and submission is independent from the program calls, therefore enabling new technologies (e.g. Hadoop) to be incorporated.

### Data security & reusability

The framework separates project specific files from reference data, scripts, and software suites that are reusable in other projects. Access to confidential data is handled transparently via the underlying Linux permission system. The transaction between projects and framework is facilitated by a project specific configuration file that defines paths to reference data as well as the analysis tasks to perform.

### Hot swapping & adaptability

NGSANE provides a unified framework (i.e. defined folder structure) for processing raw data from different experimental protocols. This allows co-investigators and reviewers to easily understand, and reproduce work using NGSANE's log and report files.

### Reproducibility & checkpoint recovery

A full audit trail is generated recording performed tasks, utilised reference data, timestamps, software version as well as HPC log files, including any errors. NGSANE gracefully recovers from unsuccessfully executed jobs be it due to failed commands, missing or incorrect input or under-resourced HPC jobs by cleanly restarting after the most recent successfully executed checkpoint.

### Robust execution & full monitoring

In our experience, resource-intensive workflows are executed in stages with interjacent human quality control, NGSANE hence focuses on providing robust checkpointing and intuitive report generation. However, workflows can be fully automated by utilising NGSANE's control over HPC-queuing systems and by leveraging the customisable interfaces between modules when submitting multiple dependent stages at once.

### Complete customisation

NGSANE's configuration file contains details about the submission system, typical resource allocations and location of third-party software. However, NGSANE's credo is that every parameter can be overwritten, hence default parameters can be adjusted in the project specific configuration file to indicate different software versions, additional resources or an altered output location.

### Automated project summary creation

NGSANE generates a highlevel summary to enable informed decisions about the experimental success. This report provides an access point for new lab members or collaborators. Furthermore, the project card can be used as gold standard for software development when employing a continuous integration server.
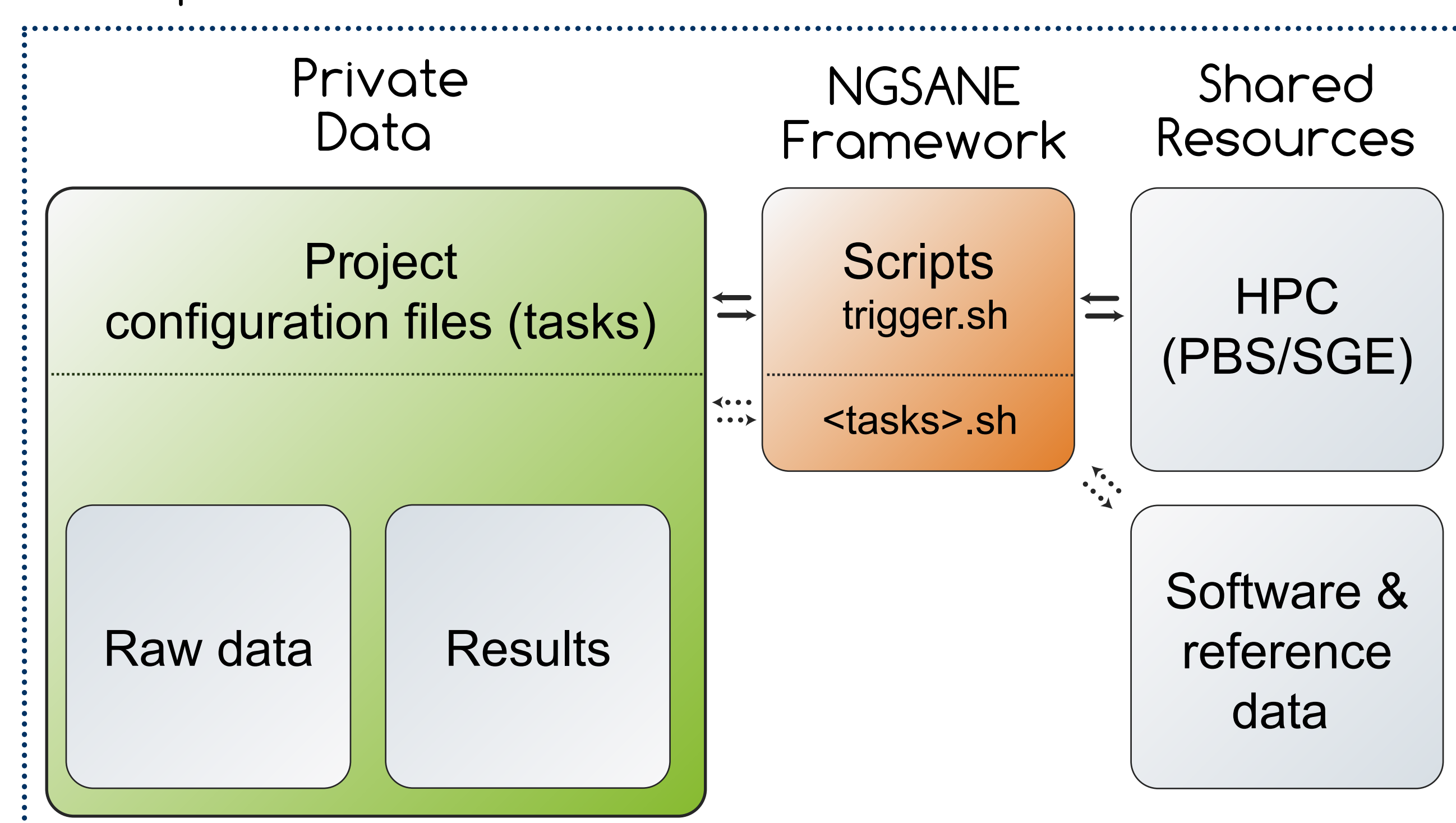
### Repeated calls

Pipelines often have to be rerun on the full or a subset of the data with possibly altered parameter settings. NGSANE facilitates and documents this by allowing multiple configuration files.
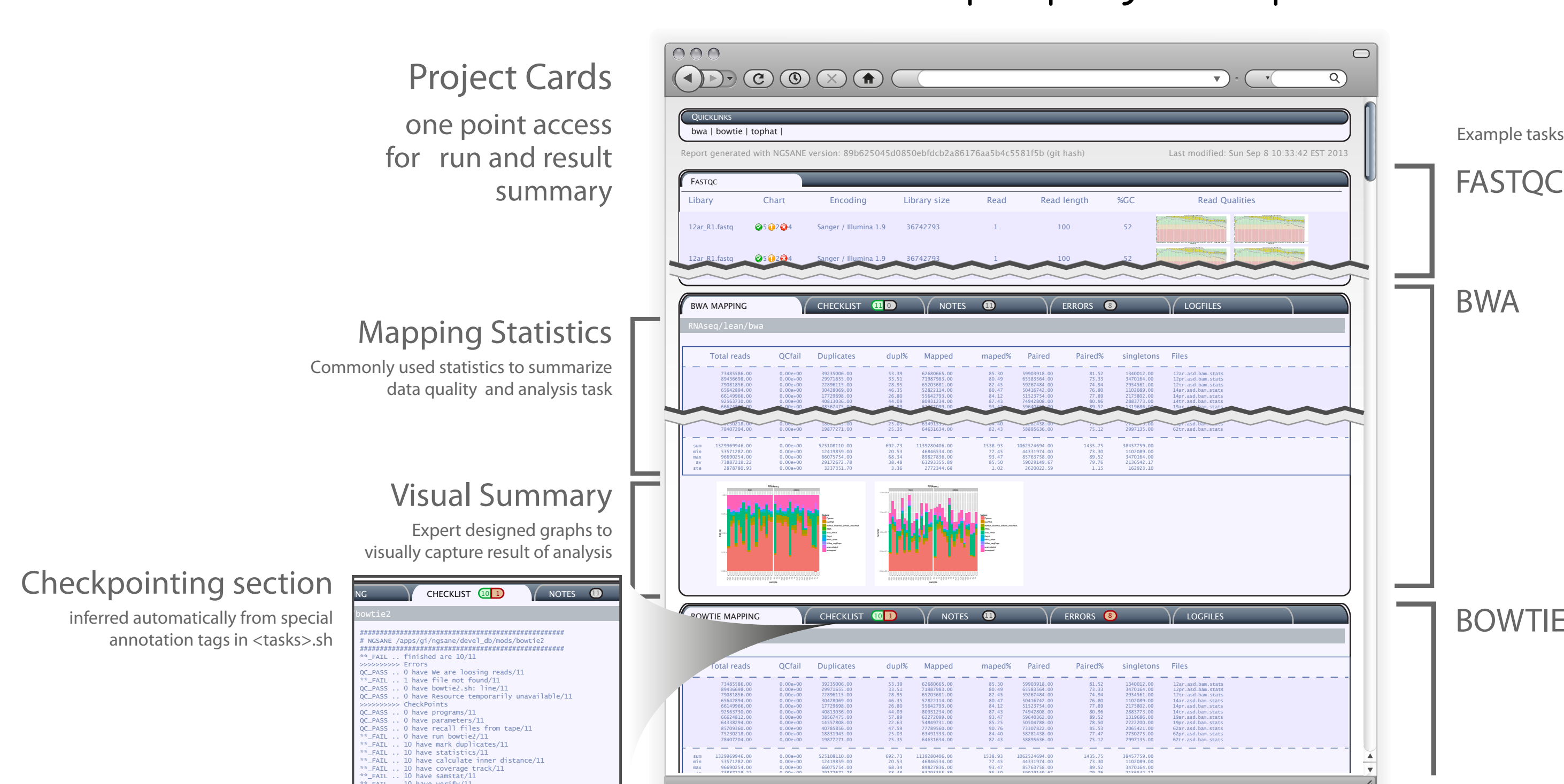
### Knowledge transfer

NGSANE provides a unified framework (i.e. defined folder structure) for processing raw data from different experimental protocols. This allows co-investigators and reviewers to easily understand, and reproduce work using NGSANE's log and report files.

## Concept



Private Data — Project configuration files (tasks) — Raw data — Results

NGSANE Framework — Scripts trigger.sh — <tasks>.sh

Shared Resources — HPC (PBS/SGE) — Software & reference data

## Example project report



Project Cards — one point access for run and result summary

Mapping Statistics — Commonly used statistics to summarize data quality and analysis task

Visual Summary — Expert designed graphs to visually capture result of analysis

Checkpointing section — inferred automatically from special annotation tags in <tasks>.sh

Example tasks — FASTQC — BWA — BOWTIE

Currently implemented workflows include those for quality control, adapter trimming, read mapping, peak calling, motif discovery, transcript assembly, variant calling and chromatin conformation analysis. These workflows utilise publicly available published software, yet allow the end user to add their own code and create new workflows as required.

NGSANE is open source and available via GitHub at https://github.com/BauerLab/ngsane

Fabian Buske
✉ f.buske@garvan.org.au
www.garvan.org.au

Denis Bauer
✉ denis.bauer@csiro.au
www.csiro.au/CCI

github SOCIAL CODING