

# Pipelines in Perl With eHive

YAPC::EU 2016

24<sup>th</sup> – 26<sup>th</sup> August 2016

Andy Yates: Team Leader, EMBL-EBI

[www.ebi.ac.uk](http://www.ebi.ac.uk)

[www.ensembl.org](http://www.ensembl.org)

<https://github.com/Ensembl/ensembl-hive>

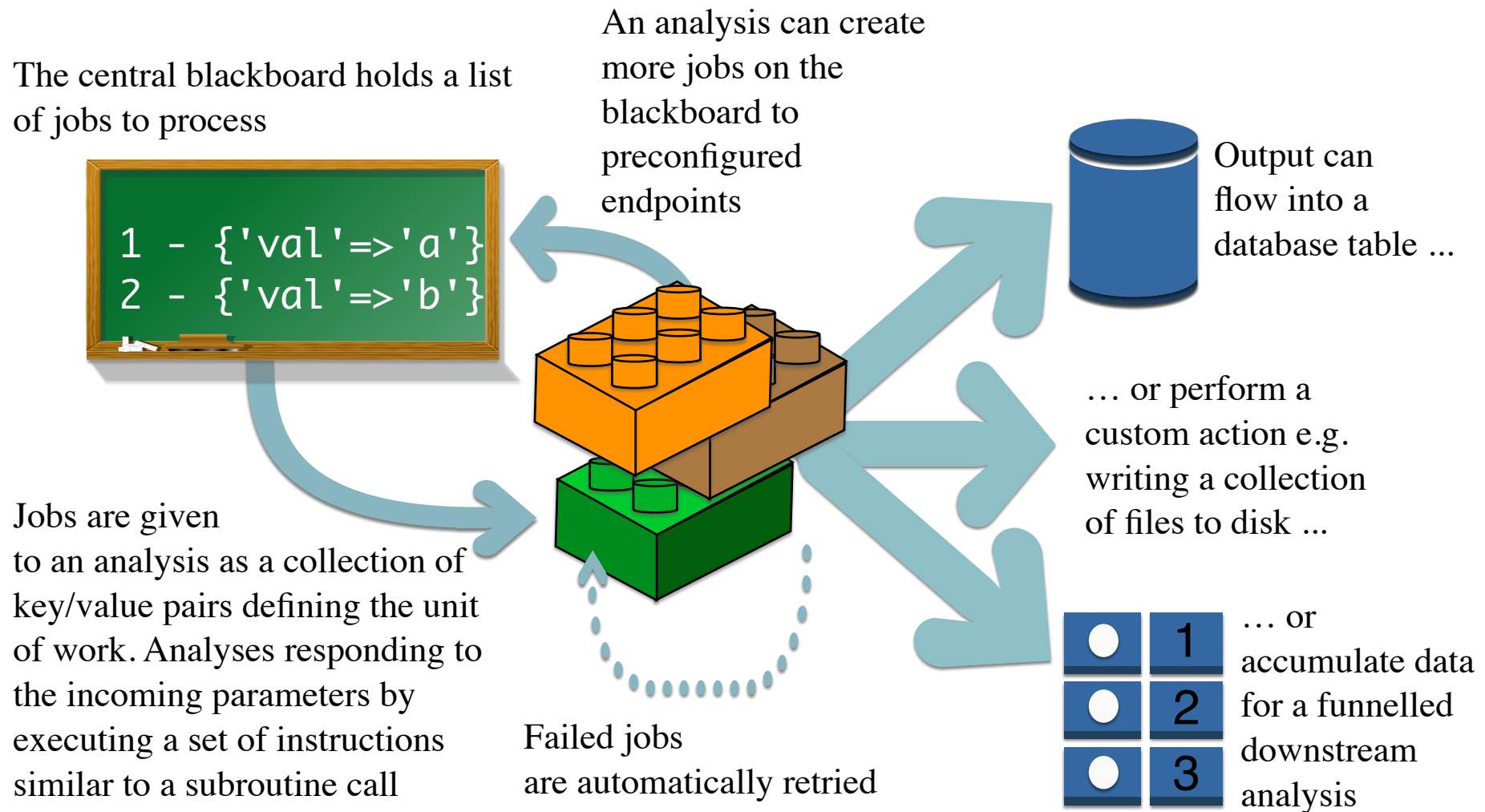
@enscore

# Processing Workflows with Perl and eHive

- <https://github.com/Ensembl/ensembl-hive>
- A workflow management system in Perl
- Uses traditional schedulers (LSF/OpenLava) to submit work
- Workflows are DAGs using semaphores to replicate programming constructs



# eHive: How it Works



```

package CountATGC;
-
use strict;
use warnings;
use Bio::SeqIO;
use base qw/Bio::Ensembl::Hive::Process/;
-
sub fetch_input {
    my ($self) = @_;
    my $chunkfile = $self->param_required('chunk_name');
    $self->param('chunk_in', Bio::SeqIO->new(-file => $chunkfile));
    return;
}
-
sub run {
    my ($self) = @_;
    my ($at, $gc) = (0, 0);
    foreach my $chunkseq ($self->param('chunk_in')->next_seq()) {
        my $seqstring = $chunkseq->seq();
        $at += @{$seqstring =~ /[AaTt]/g};
        $gc += @{$seqstring =~ /[GgCc]/g};
    }
    $self->param('at', $at);
    $self->param('gc', $gc);
    return;
}
-
sub write_output {
    my ($self) = @_;
    $self->dataflow_output_id({
        at => $self->param('at'),
        gc => $self->param('gc'),
    });
    return;
}
-
1;

```

## Get your input data sets

Get a path to a data file, initialise another object for reading and stash it

## Do your processing

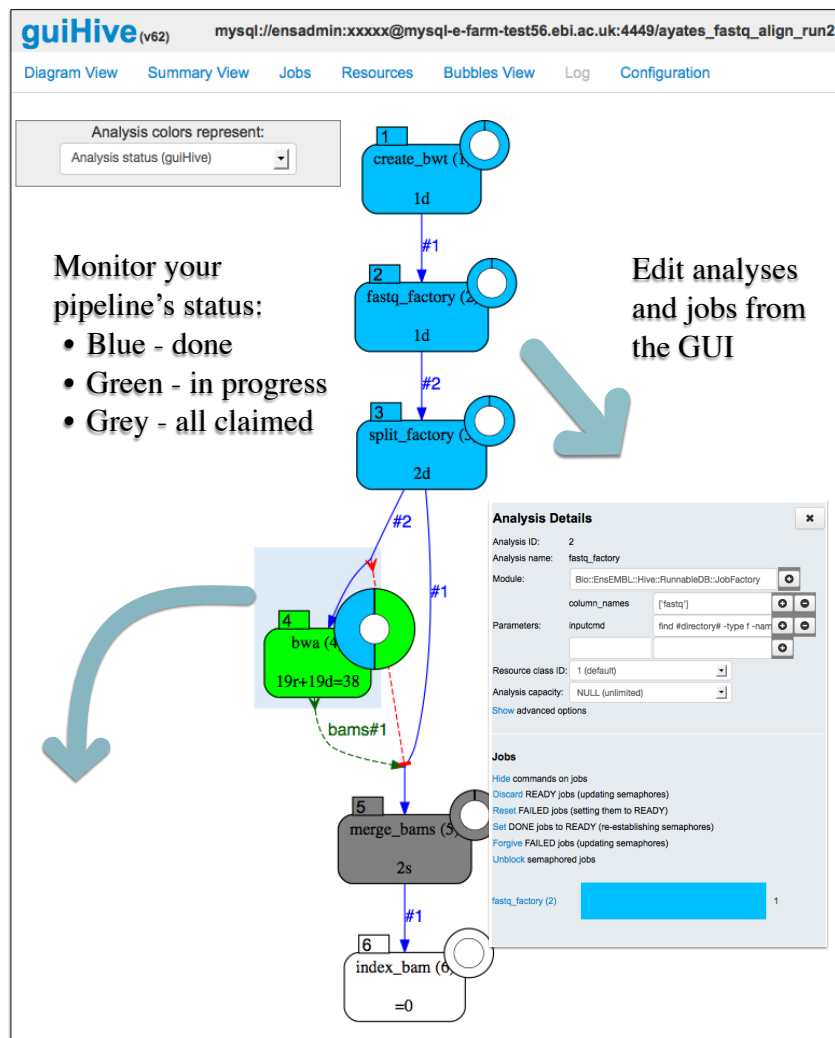
We are using the object to get sequence, counting the number of AT/GC by regex and stashing the result

## Send a hash to the next task

Current task does not need to know what will consume the hash



# guiHive for Monitoring Your Application



- Written in GO with D3.js
- Persistent web server talking to multiple hives
- Edit your pipeline as it is running
- Watch your application execute if you really want to

Tutorial from

<https://github.com/andrewyatz/eHiveDemo-yapceu2016>

# The Entire Ensembl Team

# Acknowledgements

**Andrew Yates<sup>1</sup>, Wasiu Akanni<sup>1</sup>, M. Ridwan Amode<sup>1</sup>, Daniel Barrell<sup>1,2</sup>, Konstantinos Billis<sup>1</sup>, Denise Carvalho-Silva<sup>1</sup>, Carla Cummins<sup>1</sup>, Peter Clapham<sup>2</sup>, Stephen Fitzgerald<sup>1</sup>, Laurent Gil<sup>1</sup>, Carlos García Girón<sup>1</sup>, Leo Gordon<sup>1</sup>, Thibaut Hourlier<sup>1</sup>, Sarah E. Hunt<sup>1</sup>, Sophie H. Janacek<sup>1</sup>, Nathan Johnson<sup>1</sup>, Thomas Juettemann<sup>1</sup>, Stephen Keenan<sup>1</sup>, Ilias Lavidas<sup>1</sup>, Fergal J. Martin<sup>1</sup>, Thomas Maurel<sup>1</sup>, William McLaren<sup>1</sup>, Daniel N. Murphy<sup>1</sup>, Rishi Nag<sup>1</sup>, Michael Nuhn<sup>1</sup>, Anne Parker<sup>1</sup>, Mateus Patricio<sup>1</sup>, Miguel Pignatelli<sup>1</sup>, Matthew Rahtz<sup>2</sup>, Harpreet Singh Riat<sup>1</sup>, Daniel Sheppard<sup>1</sup>, Kieron Taylor<sup>1</sup>, Anja Thormann<sup>1</sup>, Alessandro Vullo<sup>1</sup>, Steven P. Wilder<sup>1</sup>, Amonida Zadissa<sup>1</sup>, Ewan Birney<sup>1</sup>, Jennifer Harrow<sup>2</sup>, Matthieu Muffato<sup>1</sup>, Emily Perry<sup>1</sup>, Magali Ruffier<sup>1</sup>, Giulietta Spudich<sup>1</sup>, Stephen J. Trevanion<sup>1</sup>, Fiona Cunningham<sup>1</sup>, Bronwen L. Aken<sup>1</sup>, Daniel R. Zerbino<sup>1</sup> and Paul Flicek<sup>1,2,\*</sup>**

<sup>1</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK and <sup>2</sup>Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

## Funding



EMBL-EBI

