

Big Data with ADAMS

What the heck is ADAMS?

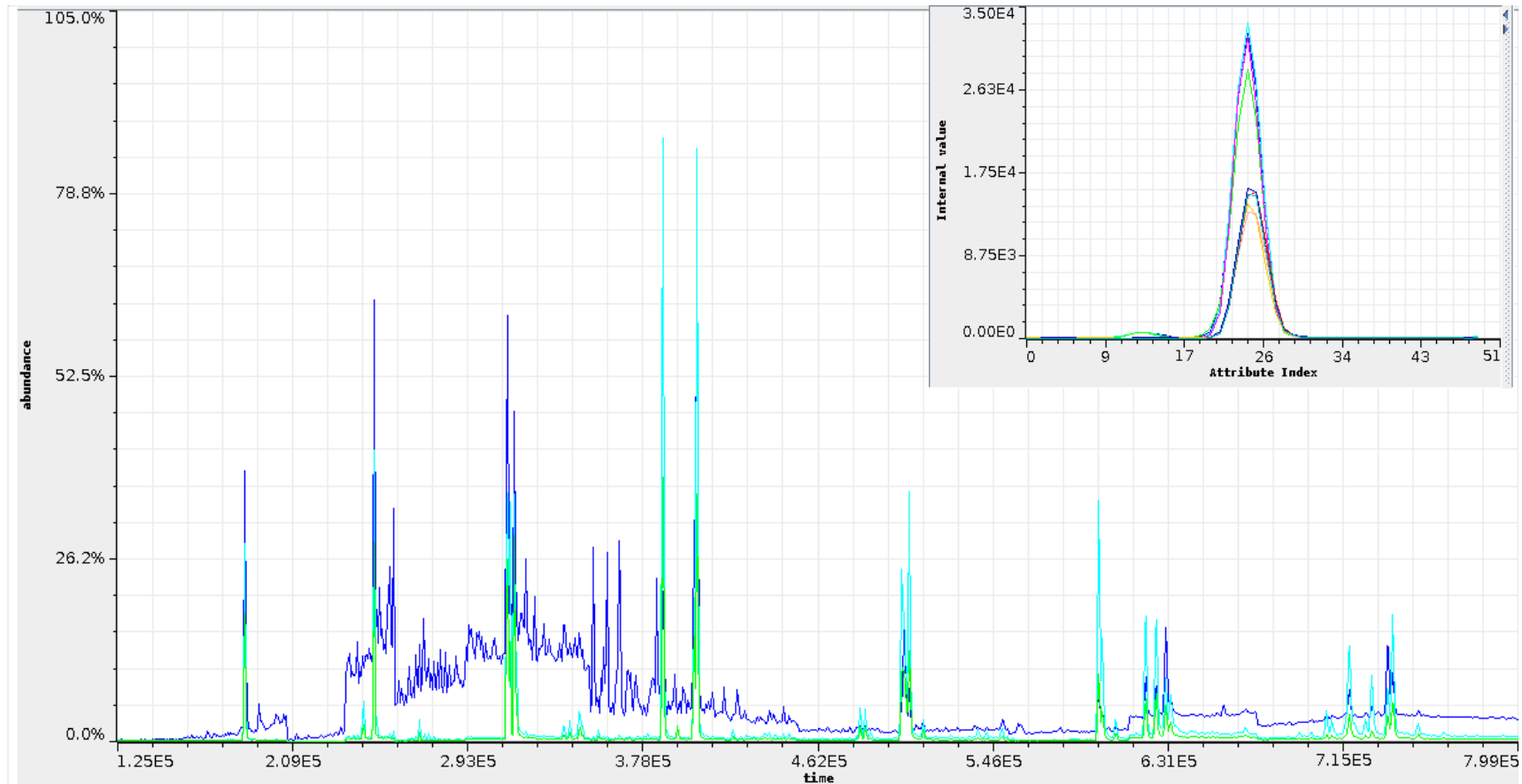
About myself

- University of Waikato, NZ - since 2003
- Senior Research Programmer
- Lead developer of **ADAMS**
- Other projects
 - **WEKA**
 - **WEKA MOOC**
 - **python-weka-wrapper**
 - **MEKA**
- More info

<http://www.cms.waikato.ac.nz/~fracpete/>

How things started

- PAH domain (poly-aromatic hydrocarbons)



GC-MS Challenges

- Multi-dimensional
(each GC point has 2-dim MS data attached)
- Retention time shifts (x-axis, linear and non-linear)
- Baseline shifts (y-axis)
- Noise (impurities in samples)
- Coelution (shoulder peaks)
- Varying sample rates (x-axis)
- Hardware problems (detector saturation)
- Multiple compounds to predict



Solution

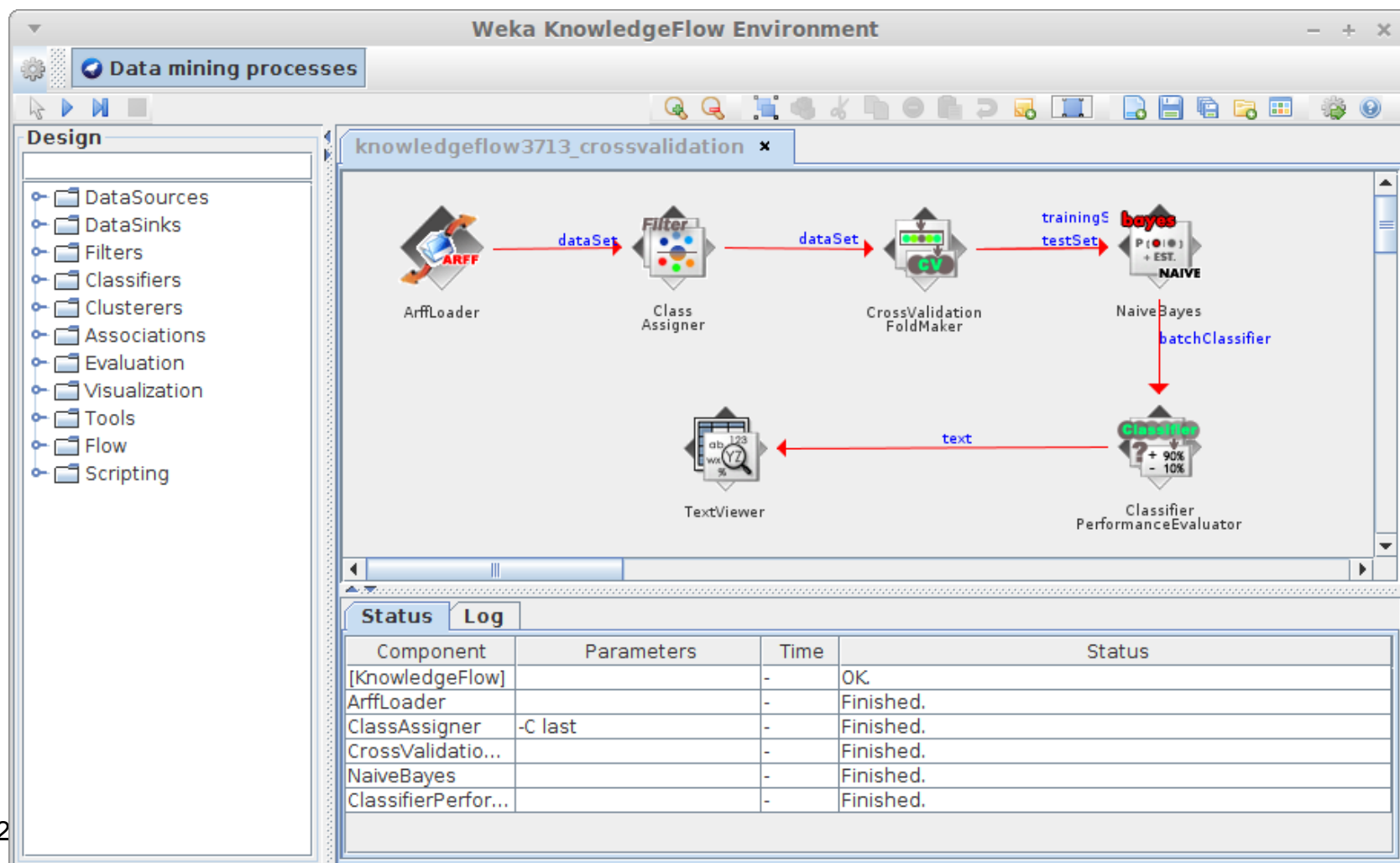
- Extensive pre-processing
- Pre-processing/prediction in parallel for various compounds
→ Workflow engine, but which one?

Workflow systems

- most workflow systems are “canvas-based”
 - user places operators on canvas
 - manually connects operators
- older systems were DAG-based
 - eg using a tree-layout
- following examples
 - cross-validation of classifier on single dataset
 - output of results

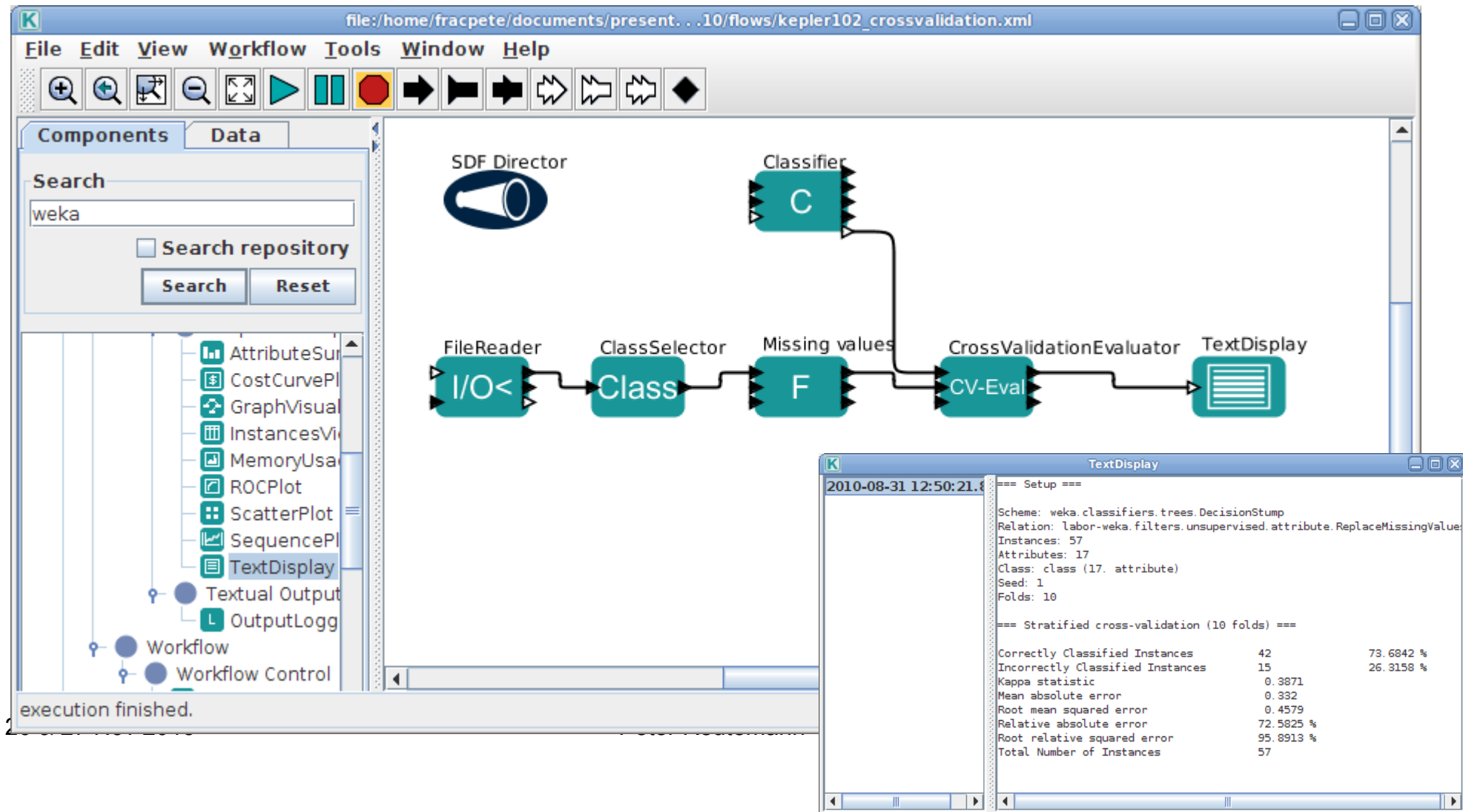
Workflow systems (2)

- WEKA KnowledgeFlow 3.7.13



Workflow systems (3)

- Kepler 1.0.2



- 26 & 27 N

04_XValidation_Nominal - RapidMiner@julia

File Edit Process Tools View Help

Result Overview PerformanceVector (ClassificationPerformance)

Table / Plot View Text View Annotations

Criterion Selector

classification_error

Multiclass Classification Performance Annotations

Table View Plot View

classification_error: 7.50% +/- 10.00% (mikro: 7.50%)

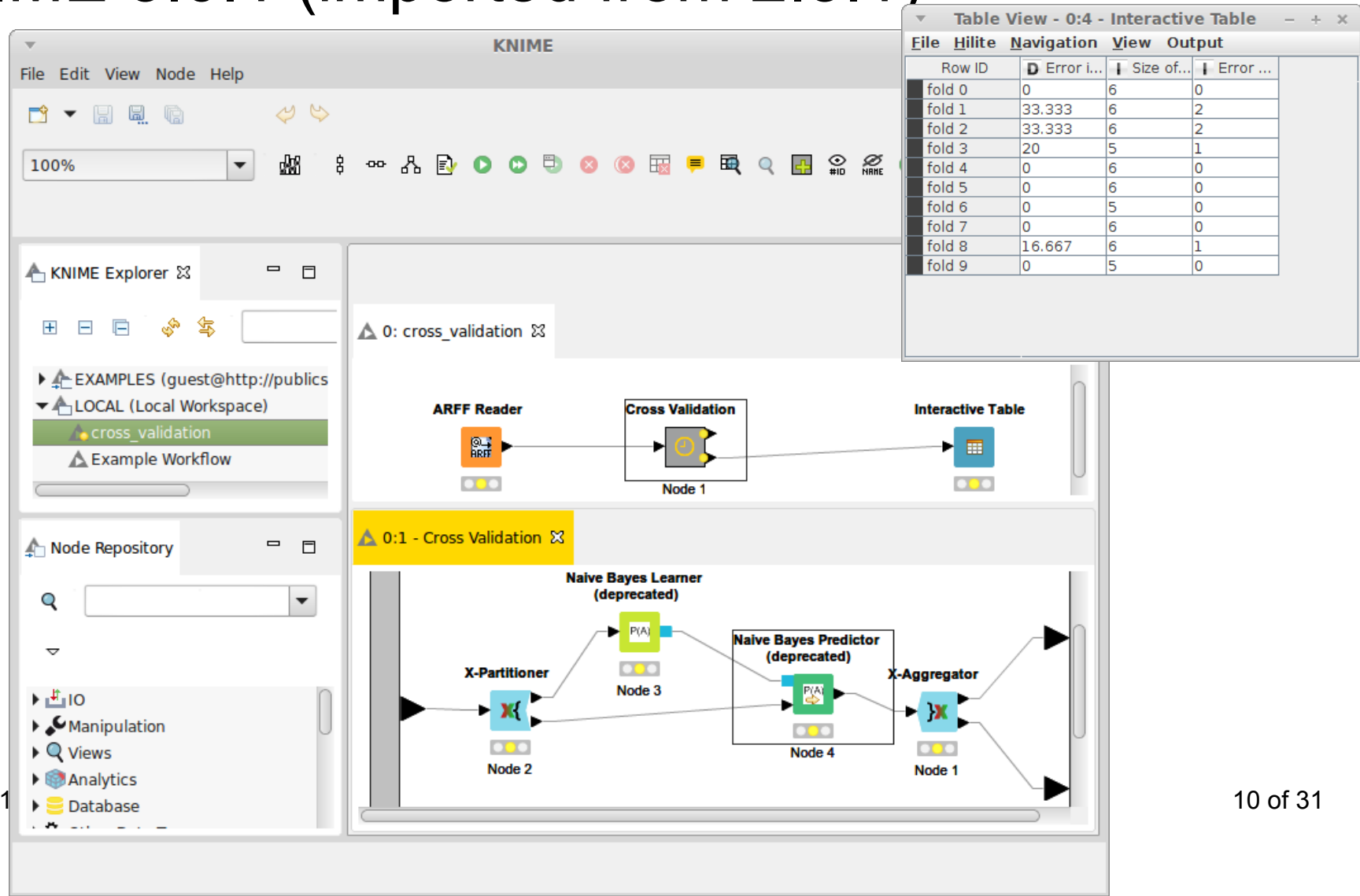
	true bad	true good
pred. bad	12	1
pred. good	2	25
class recall	85.71%	96.15%

Log

Aug 31, 2010 11:24:37 AM WARNING Error checking for updates: javax.xml.ws.WebServiceException: Failed to access

Workflow systems (5)

- KNIME 3.0.1 (imported from 2.3.1)



The screenshot displays the KNIME 3.0.1 interface. The main workspace shows a workflow with the following nodes:

- ARFF Reader** (Node 1)
- Cross Validation** (Node 1)
- Interactive Table** (Node 1)
- X-Partitioner** (Node 2)
- Naive Bayes Learner (deprecated)** (Node 3)
- Naive Bayes Predictor (deprecated)** (Node 4)
- X-Aggregator** (Node 1)

The **Table View - 0:4 - Interactive Table** window is open, showing the following data:

Row ID	Error i...	Size of...	Error ...
fold 0	0	6	0
fold 1	33.333	6	2
fold 2	33.333	6	2
fold 3	20	5	1
fold 4	0	6	0
fold 5	0	6	0
fold 6	0	5	0
fold 7	0	6	0
fold 8	16.667	6	1
fold 9	0	5	0

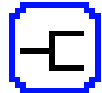







Why another workflow?

- KnowledgeFlow too WEKA-centric
- Initial development with Kepler
- Canvas-based set up is tedious
- Only minimal subset of functionality necessary (e.g., no grid computing)
- Primary connection types: 1-to-1 and 1-to-n
→ prototype implementation

What is ADAMS?

- Java, GPLv3
- Base modules
 - access, core, compress, event, excel, gnuplot, groovy, imaging (+boofcv, imagej, imagemagick, openimaj), jython, latex, maps, meta, moa, net, odf, osm, pdf, r, random, spreadsheet, timeseries, twitter, visualstats, weka
- Add-ons modules
 - heatmap, image-webservice, jooq, meka, rats, video, webservice, weka-webservice
- Incubator modules
 - nlp, jclouds, openstack, openml, javacv, ...

Flow

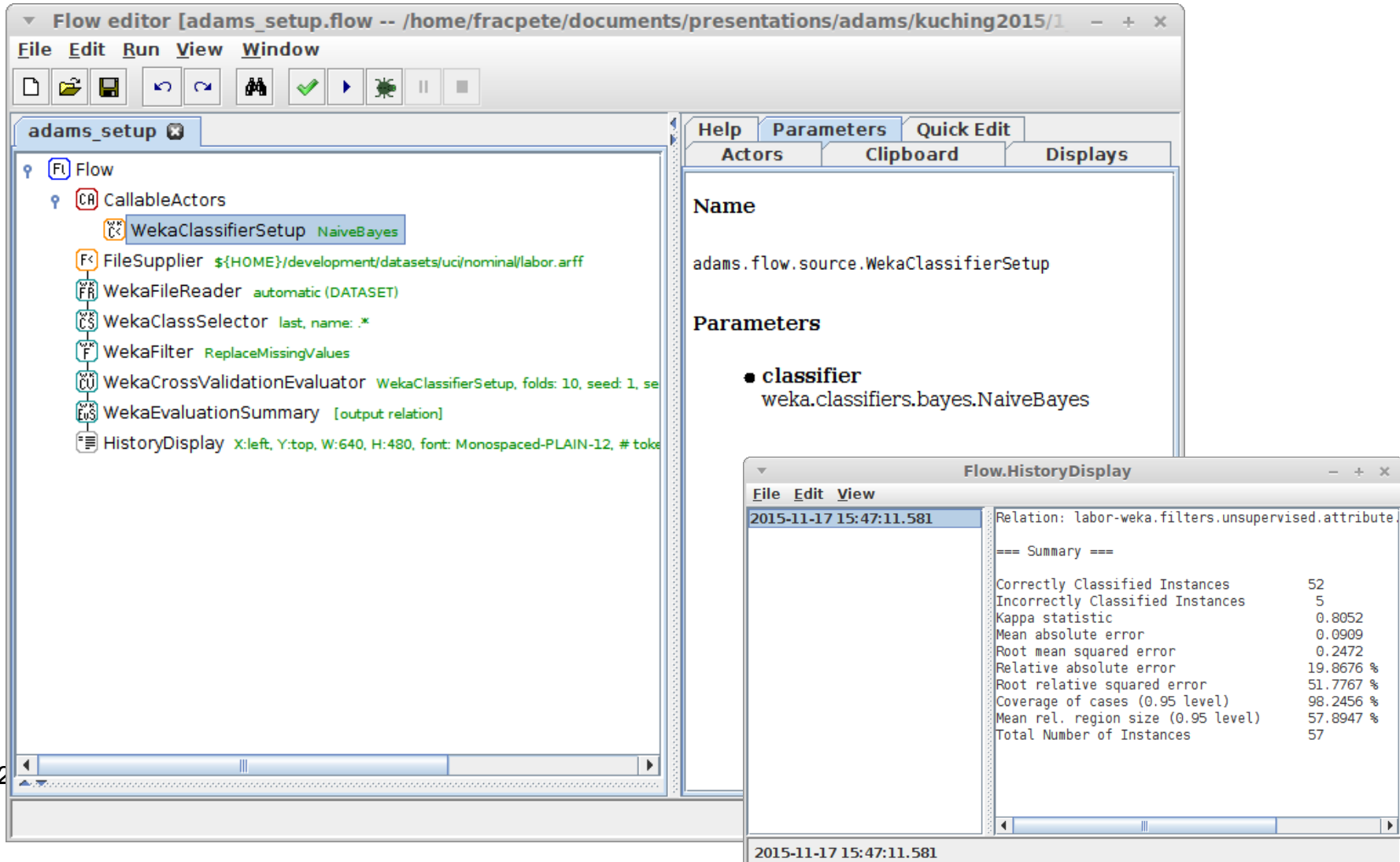
- Operators are called “actors”
- Actors arranged in tree, no connections
- Actor “handlers” nest other actors
 - e.g., sequence of actors
- Control actors control data flow
 - e.g., branch , tee , if-then-else , switch 
- Input/output defines
 - standalone , source , transformer , sink 

Flow (2)

- Data-driven system, but events possible
- Tree only supports 1-to-n connections
- Simulating n-to-m semantics
 - Containers
 - Variables
 - Internal storage
 - Callable actors

Previous example

- Same workflow in ADAMS



The screenshot shows the ADAMS Flow editor interface. The main window displays a flow graph for 'adams_setup'. The flow starts with a 'Flow' actor, followed by 'CallableActors' containing 'WekaClassifierSetup NaiveBayes'. This is followed by a sequence of actors: 'FileSupplier' (path: \${HOME}/development/datasets/uci/nominal/labor.arff), 'WekaFileReader' (automatic (DATASET)), 'WekaClassSelector' (last, name: *), 'WekaFilter' (ReplaceMissingValues), 'WekaCrossValidationEvaluator' (WekaClassifierSetup, folds: 10, seed: 1, se), and 'WekaEvaluationSummary' (output relation). The flow ends with a 'HistoryDisplay' actor (X:left, Y:top, W:640, H:480, font: Monospaced-PLAIN-12, # tokens).

The 'Parameters' panel on the right shows the 'Name' as 'adams.flow.source.WekaClassifierSetup' and the 'Parameters' as:

- classifier
weka.classifiers.bayes.NaiveBayes

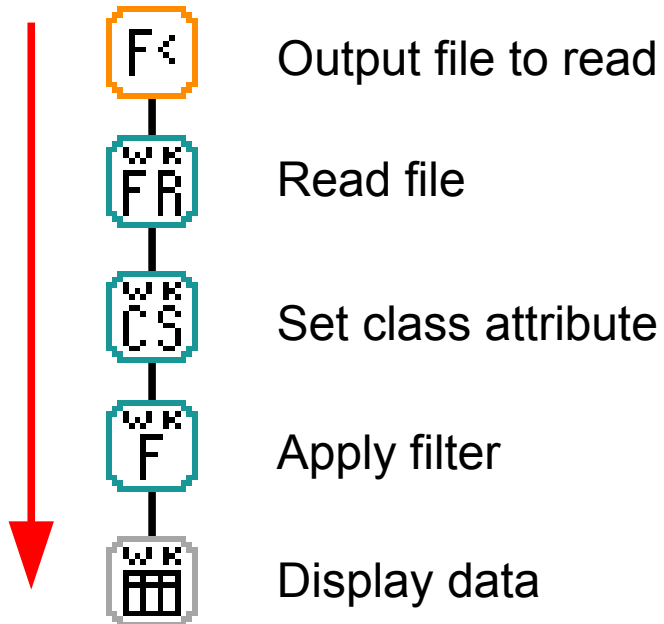
The 'Flow.HistoryDisplay' window shows the results of the evaluation. The summary is as follows:

Summary	
Correctly Classified Instances	52
Incorrectly Classified Instances	5
Kappa statistic	0.8052
Mean absolute error	0.0909
Root mean squared error	0.2472
Relative absolute error	19.8676 %
Root relative squared error	51.7767 %
Coverage of cases (0.95 level)	98.2456 %
Mean rel. region size (0.95 level)	57.8947 %
Total Number of Instances	57

Examples

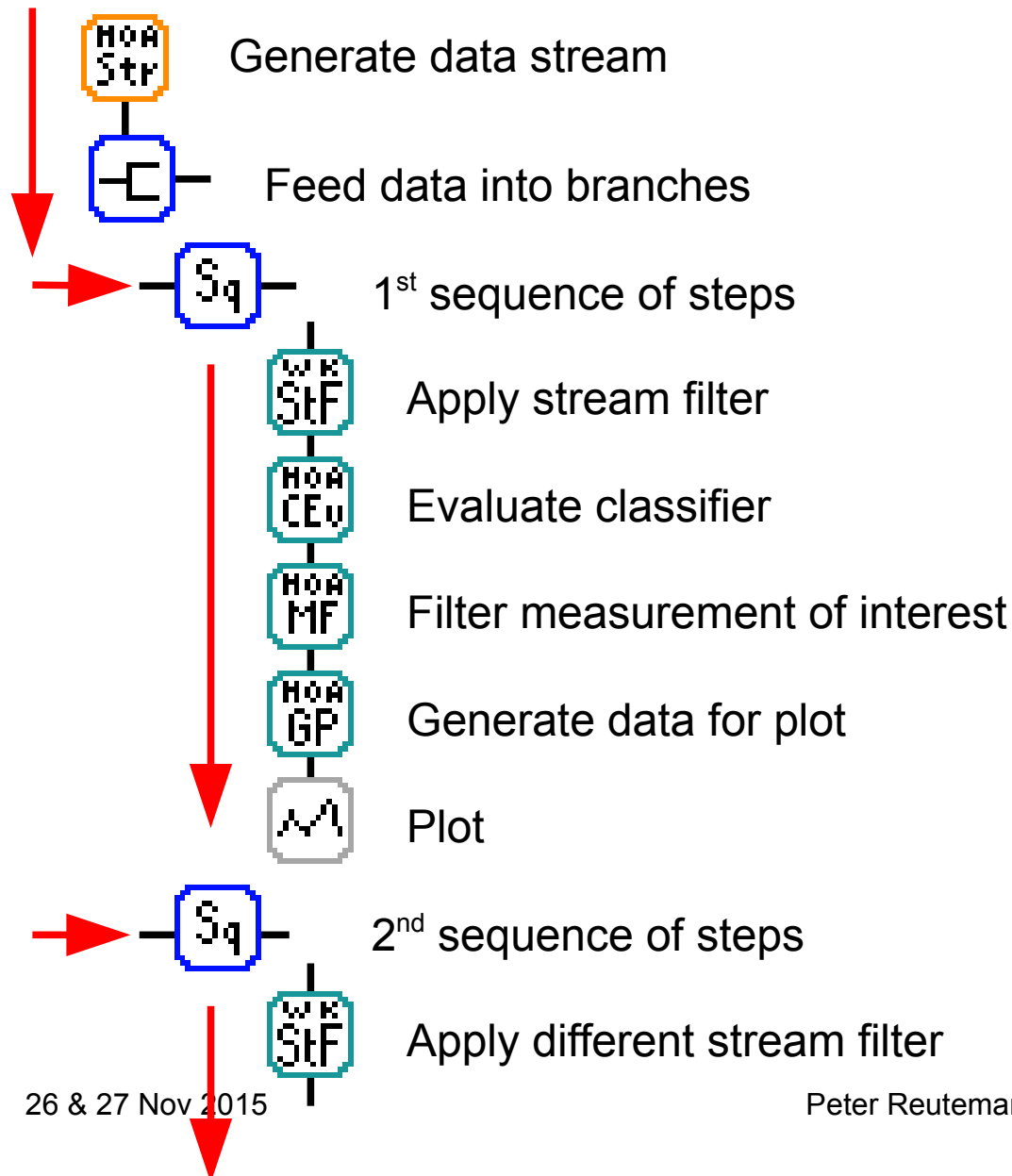


Execute nested actors one after the other



Load dataset,
apply filter and
display dataset

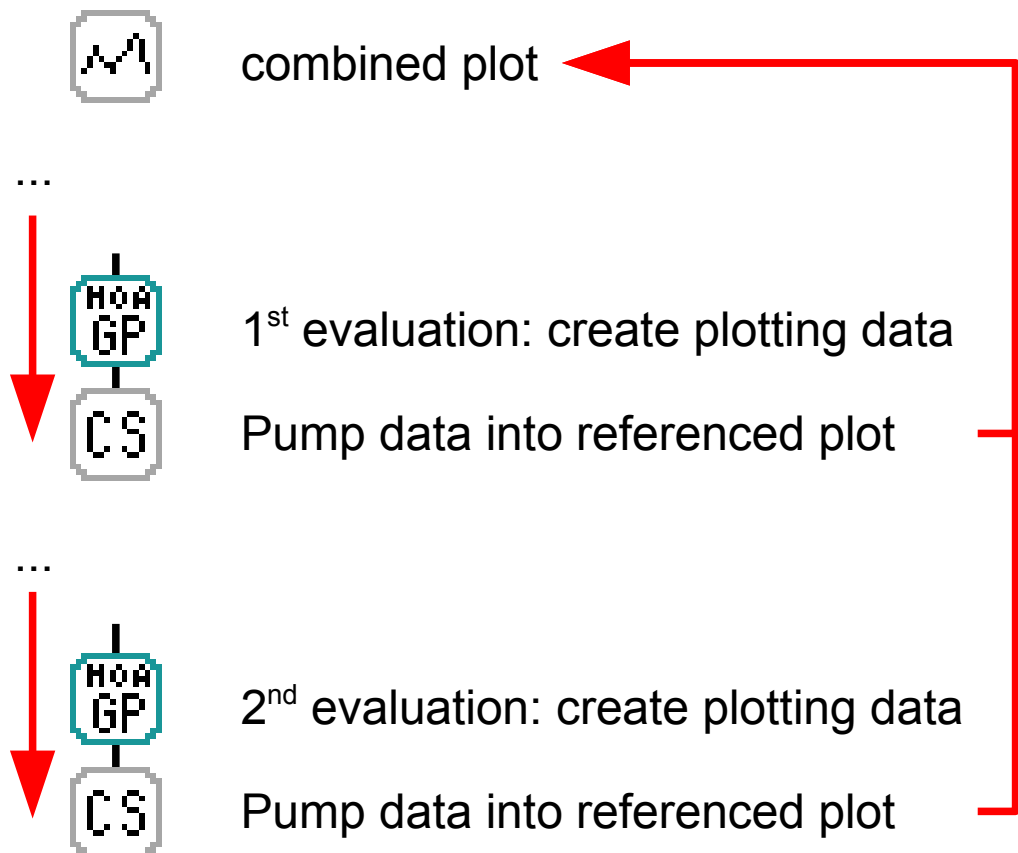
Examples (2)



Filter data stream in two separate branches with different filters, evaluate classifier and plot metric

Examples (3)

CA groups actors accessible via their name (“callable actors”)



Generate combined plot of two evaluations by using “callable actors” functionality

Tensorflow

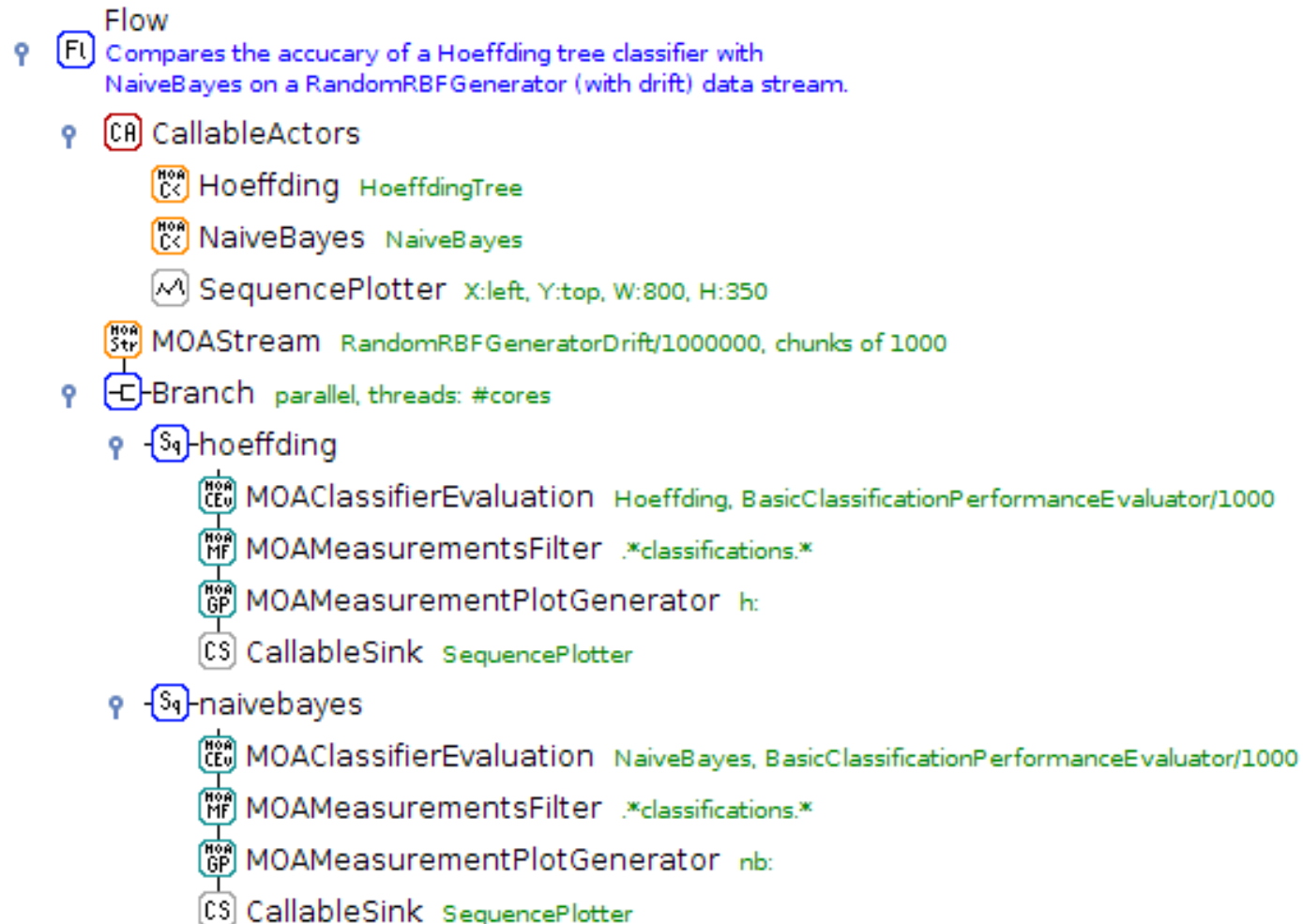
Feature	ADAMS	Tensorflow
Bindings	Java	Python and C++
GUI	✓	
Scripting	Groovy/Jython	Python
GPU		✓
Deep learning		✓
Non-deep learning	✓	
Multi-processor	✓	✓
Parallel execution	✓	✓
Remote execution	✓	
Data flow programming	kind of	✓
Platforms	Linux, Mac, Windows	Linux, Mac
License	GPLv3	Apache 2.0

TensorFlow Disappoints – Google Deep Learning falls shallow [source: KDNuggets]

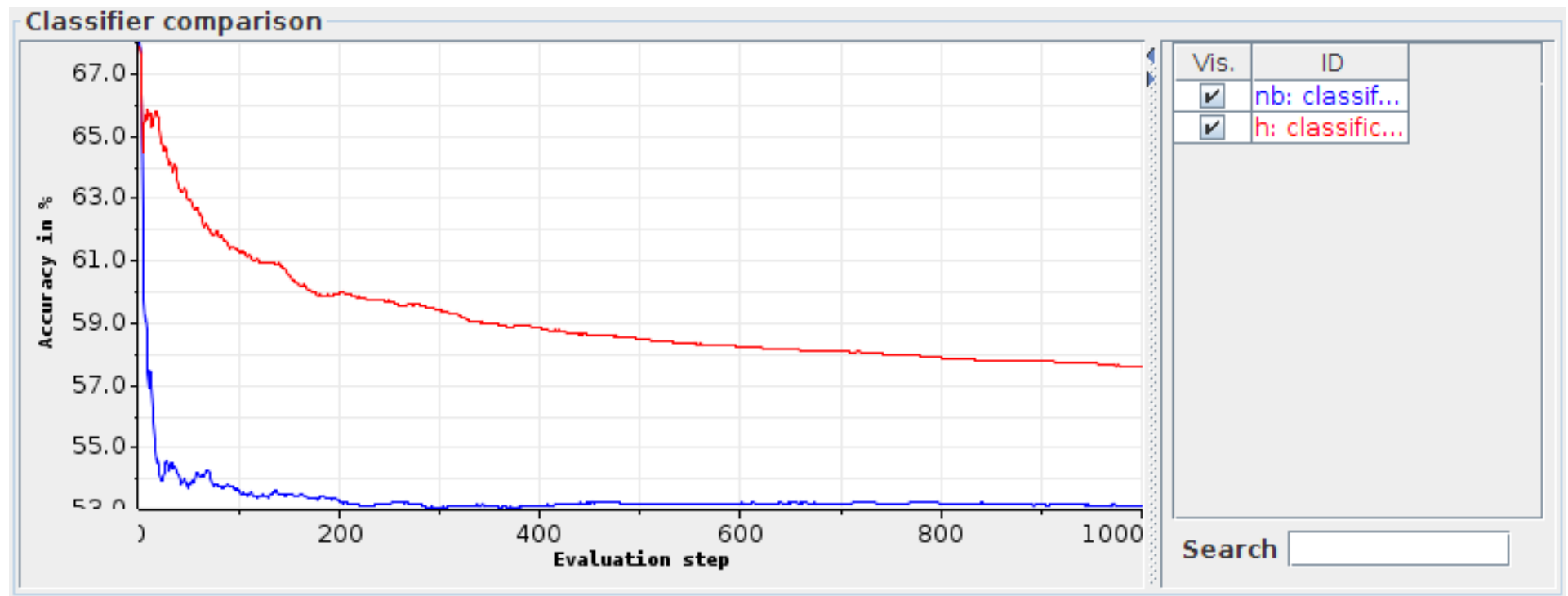
Research (demos)

- Compare two MOA classifiers (drift)
- Compare MOA classifier on different streams
- MOA cluster visualization
- Track mouse in video

MOA - Drift

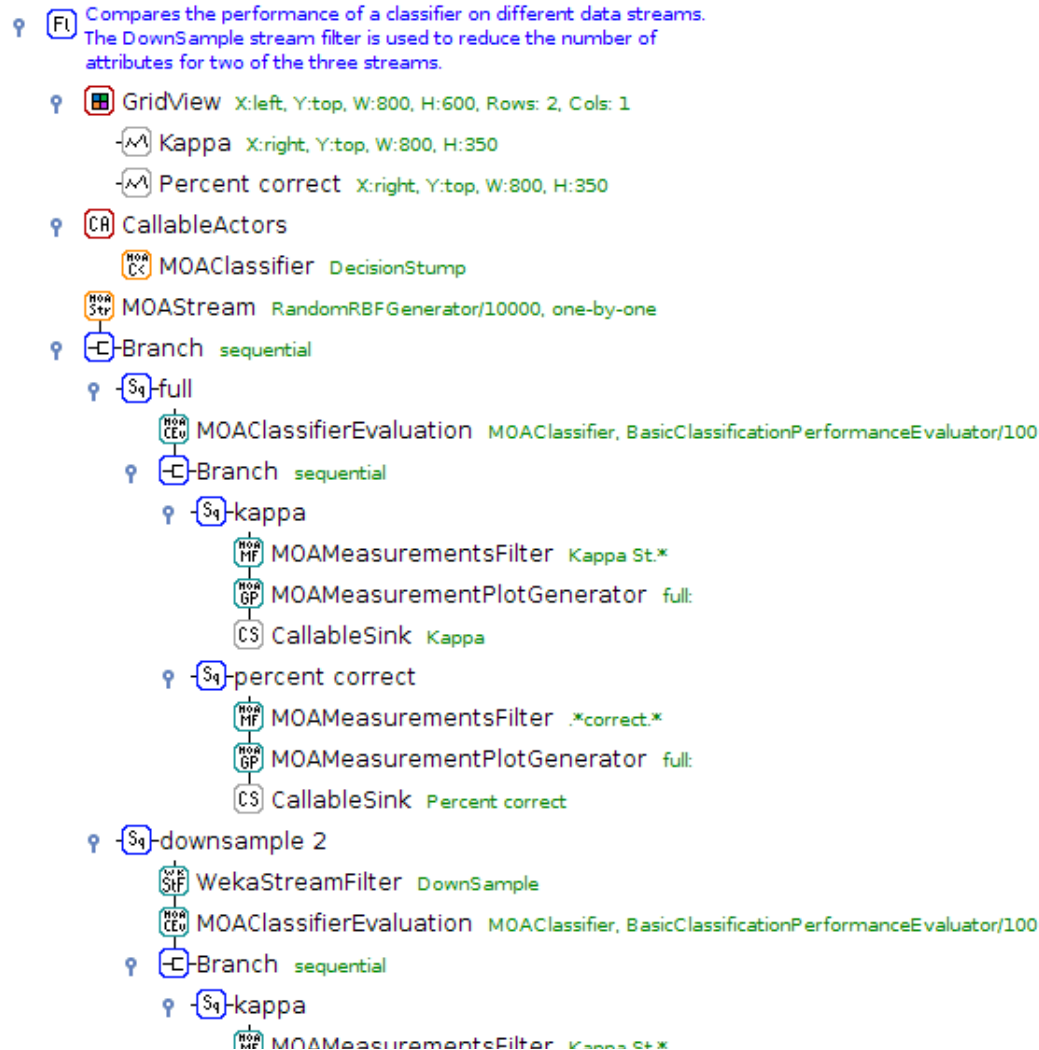


MOA - Drift

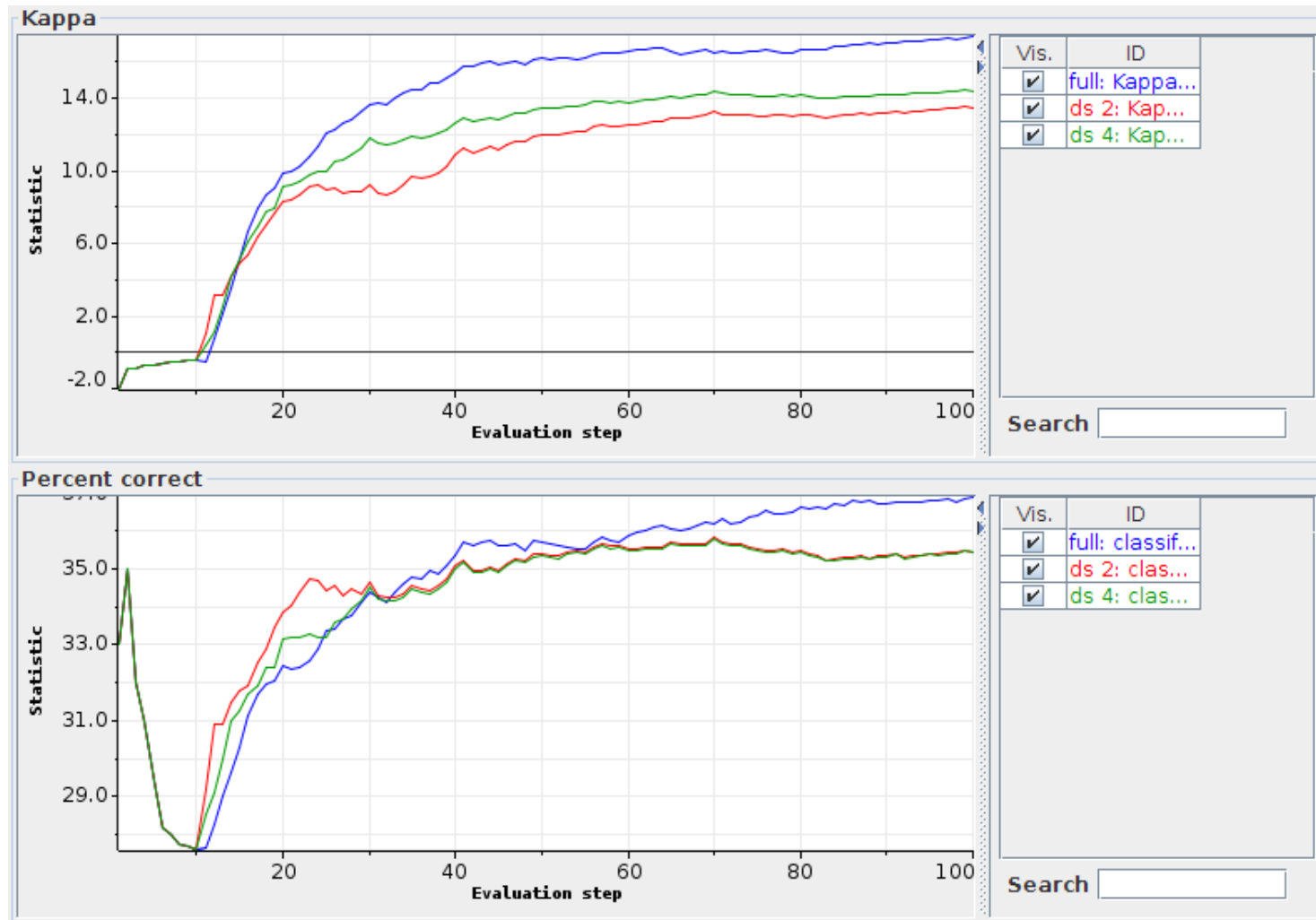


MOA - different streams

Flow

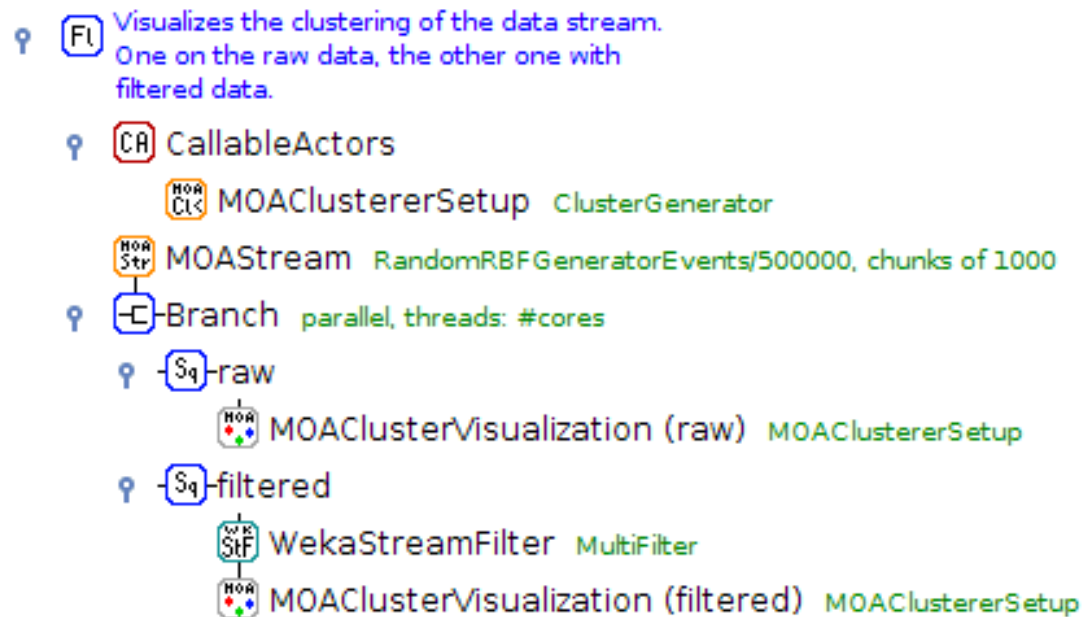


MOA - different streams

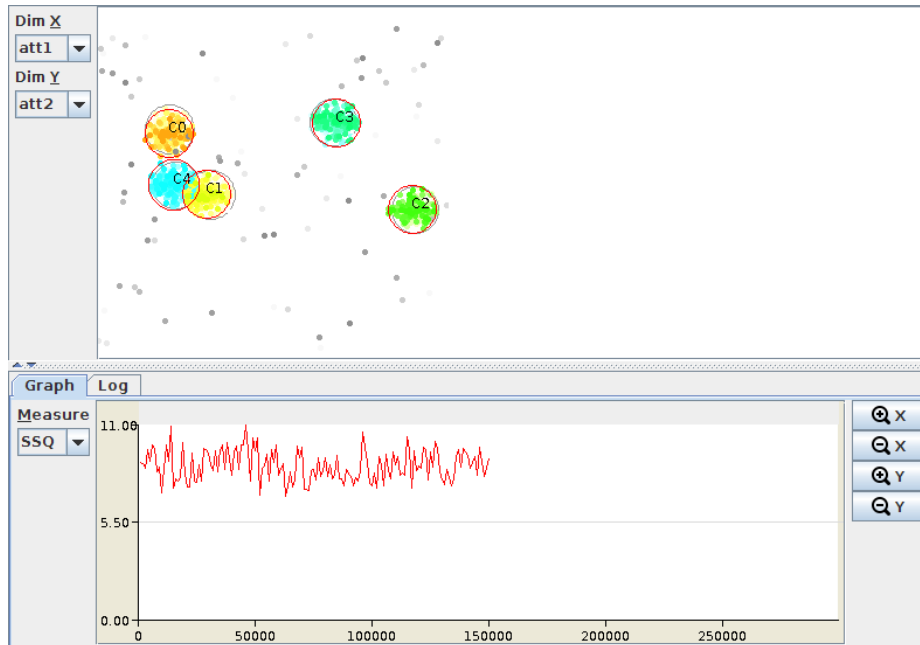


MOA - Cluster visualization

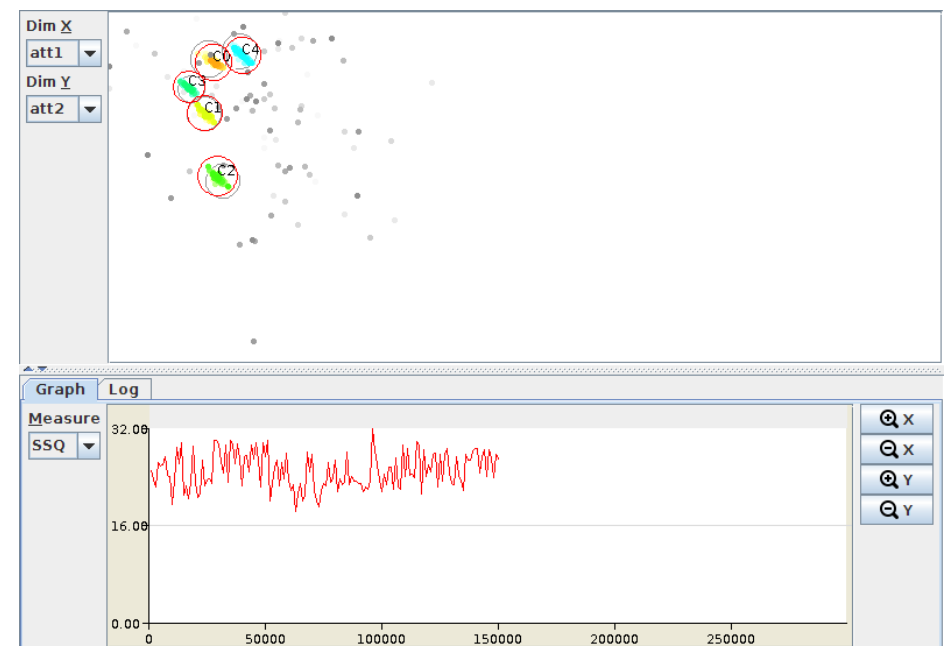
Flow



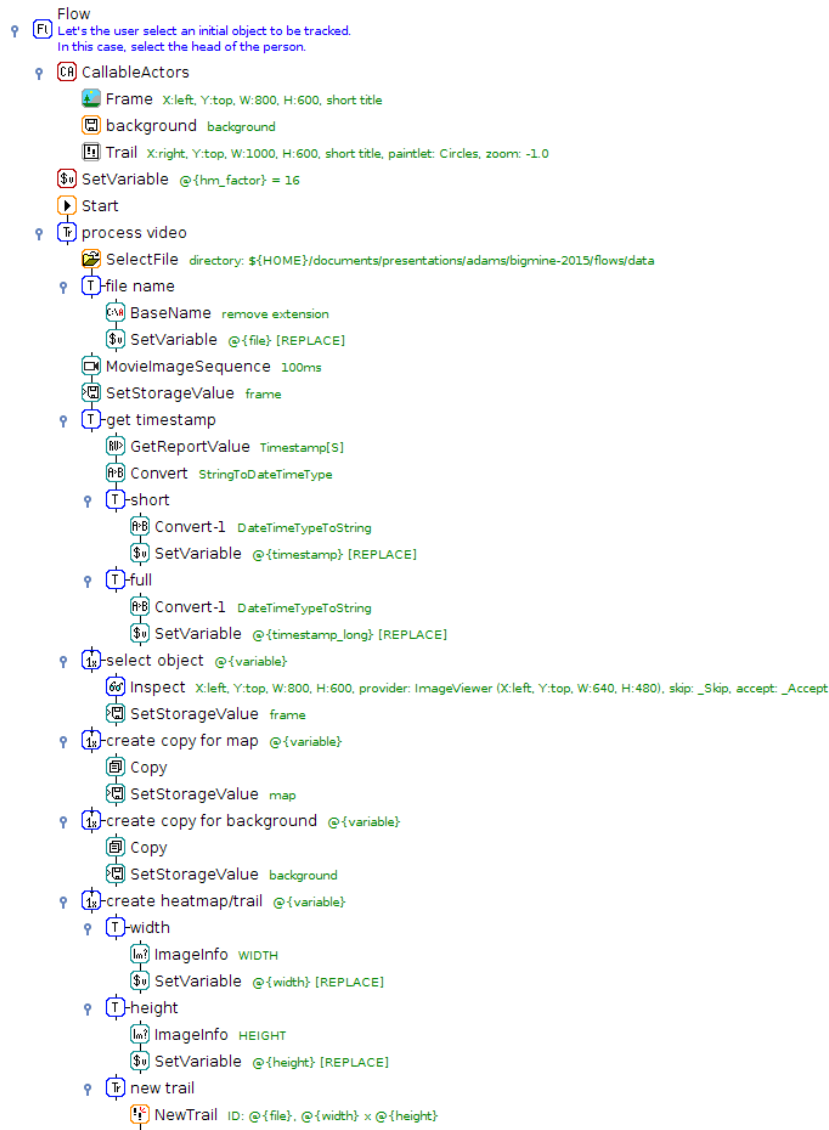
MOA - Cluster visualization



Stream 2

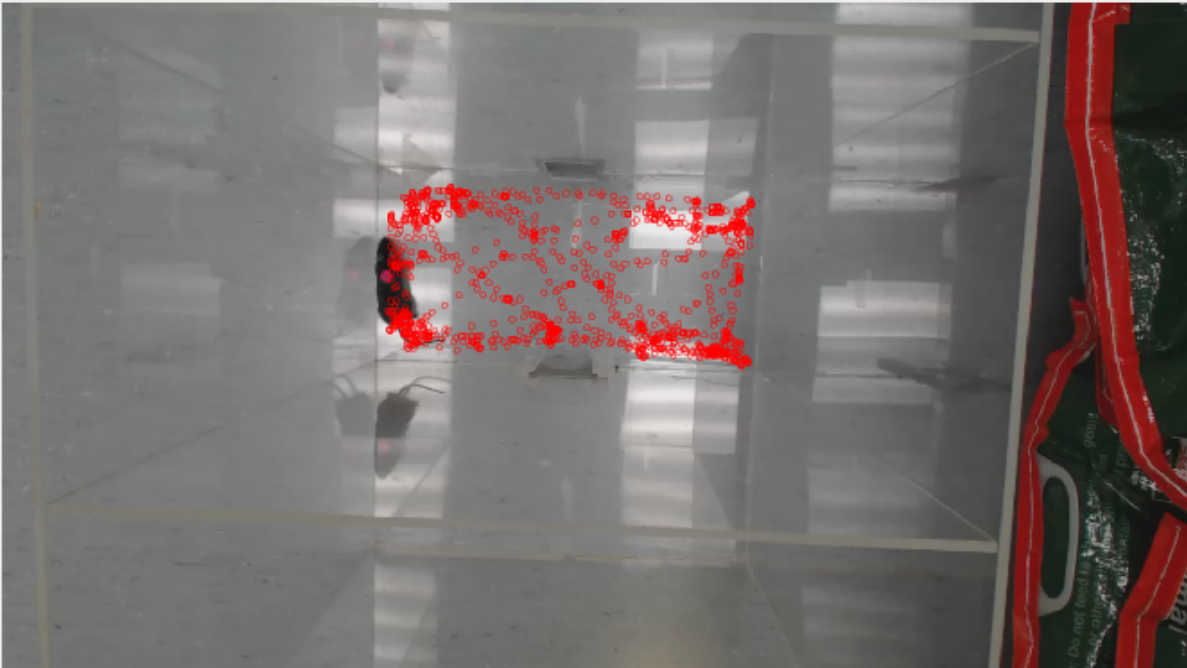


Track mouse



Track mouse

ViewTrail



X: 474 Y: 22 Zoom: 51.7%

DataLog

Name	Type	Value
Trail.Height	N	720.0
Trail.Width	N	1280.0

Search

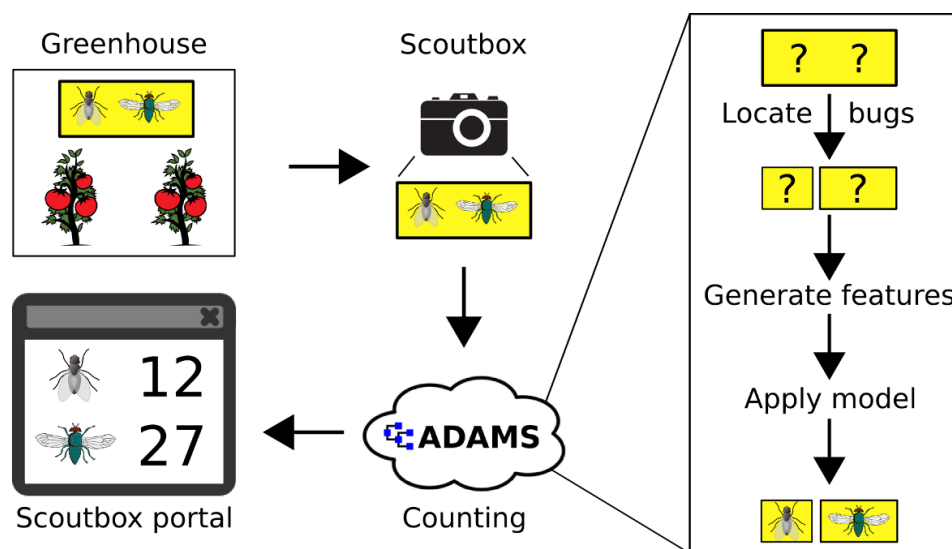
Industry

- **BLGG** - environmental lab in NL
- Spectral analysis
 - XRF: 10,000, MIR: 2,000, NIR: 1,500
- In operation since 2006
- Predictive modelling: soil, plant (~250 models)
- 1,000 to 3,000 samples per day
- Savings due to less wet chemistry
 - USD 18 million to USD 33 million per year

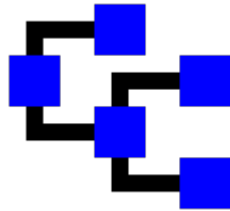
Industry (2)

- **Cropwatch BV**

- monitoring trends in insect populations (whitefly, macrolophus, thrips)
- currently used in greenhouses
- analyzing images of sticky plates



Questions?



<https://adams.cms.waikato.ac.nz/>

@TheAdamsFlow