

Workflow Mining:

Identification of frequent patterns in a large collection of KNIME workflows

Nils Weskamp, Research Scientist
Computational Chemistry
nils.weskamp@boehringer-ingenelheim.com

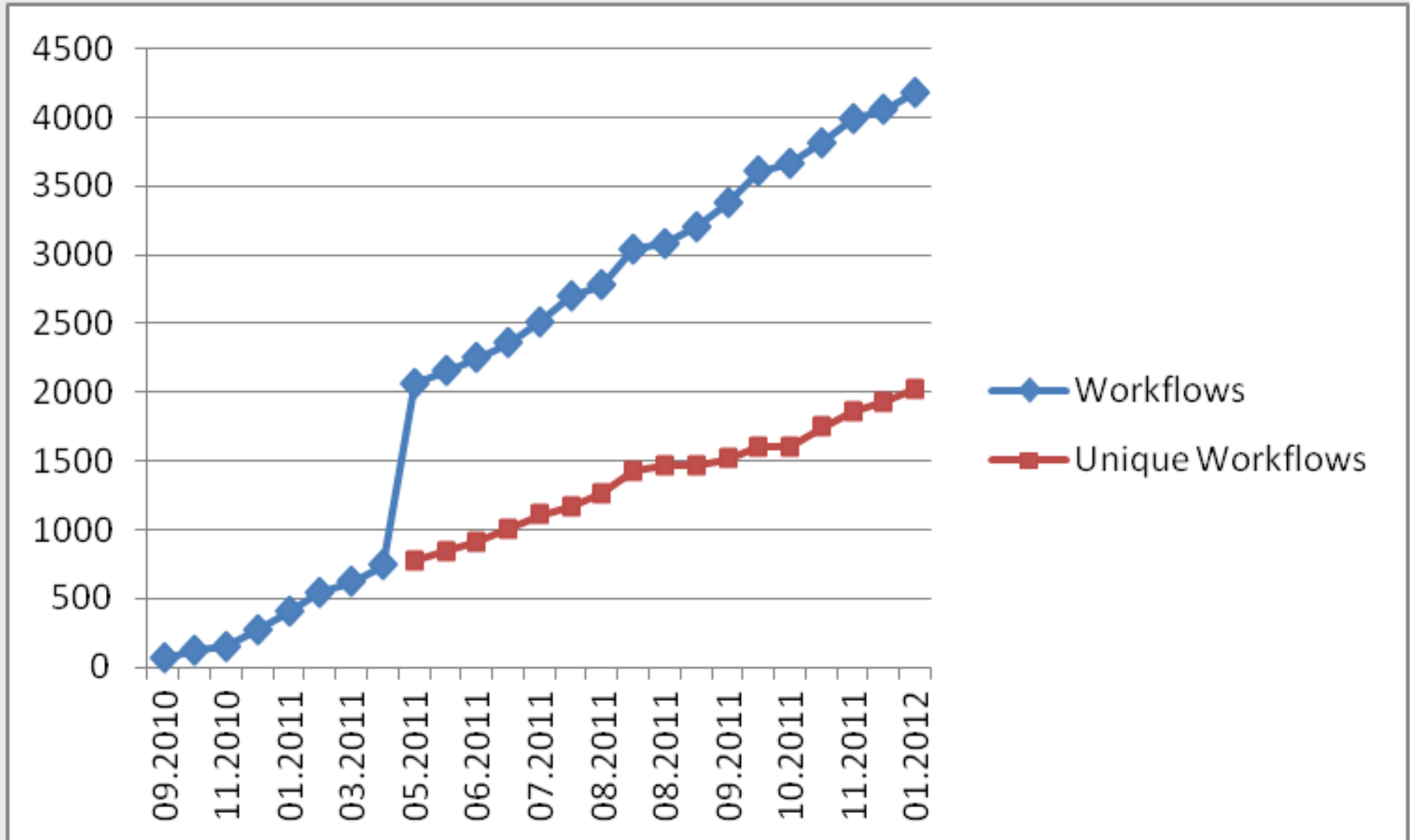


- Motivation
 - History & current status
 - Why “Workflow Mining”?
- “Static Workflow Mining”
 - Extraction of workflow graphs
 - Frequent sub-graph mining
 - Results
- “Dynamic Workflow Mining”
 - KNIME Profiler Plugin
 - Analysis of profiling logs
- Conclusions and Outlook

- **Motivation**
 - History & current status
 - Why “Workflow Mining”?
- “Static Workflow Mining”
 - Extraction of workflow graphs
 - Frequent sub-graph mining
 - Results
- “Dynamic Workflow Mining”
 - KNIME Profiler Plugin
 - Analysis of profiling logs
- Conclusions and Outlook

- 2009
 - First look at KNIME, decision to perform in-depth evaluation
 - Identification of critical functional gaps
 - Inhouse developer training
 - Start of collaboration with KNIME.com
- 2010
 - Implementation of missing nodes, integration of various external chemistry tools
 - Licensing of additional node collections from commercial vendors
 - Migration of deployed workflows (e.g. QSAR-models for ADME parameters)
- 2011
 - Migration of end-users from computational chemistry
 - Inhouse end-user trainings
 - Extension of user base to other groups (e.g. computational biology, IS, HTS)

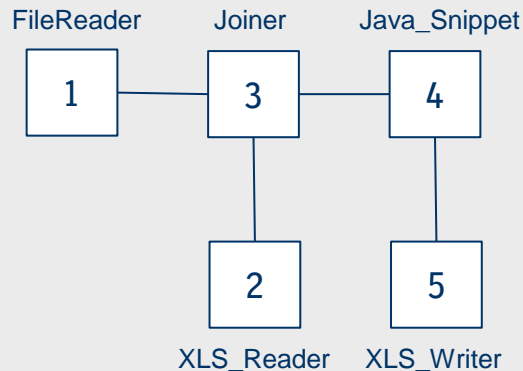
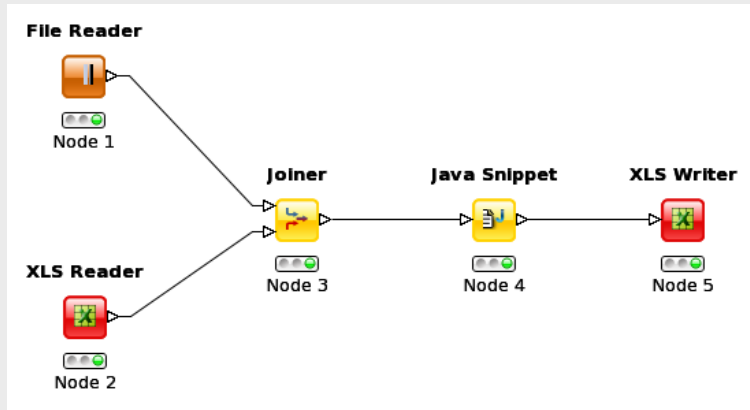
Development of KNIME usage



- A significant group of end-users is applying KNIME on a regular / daily basis
 - Consistently more than 100 new, unique workflows are created per month
 - Significant exchange and sharing of workflows among users
- A lot of feedback from end-users
 - Bug reports
 - Feature requests
 - How do I ...
 - ...
- Collecting feedback is simple, but ...
 - ... how to prioritize feature requests ?
 - ... what are best practices / common pitfalls / bottlenecks ?

- Difficult to assess significance of requests in a growing and diverse user community
 - Many different application scenarios
 - How are our users actually using KNIME?
 - Different people solve the same problem in a different way!
 - Which nodes are really important?
- Identify key users to cover diversity in user community
 - Still, a lot of gut feeling involved
- Analyze existing workflow collection to get reliable usage statistics

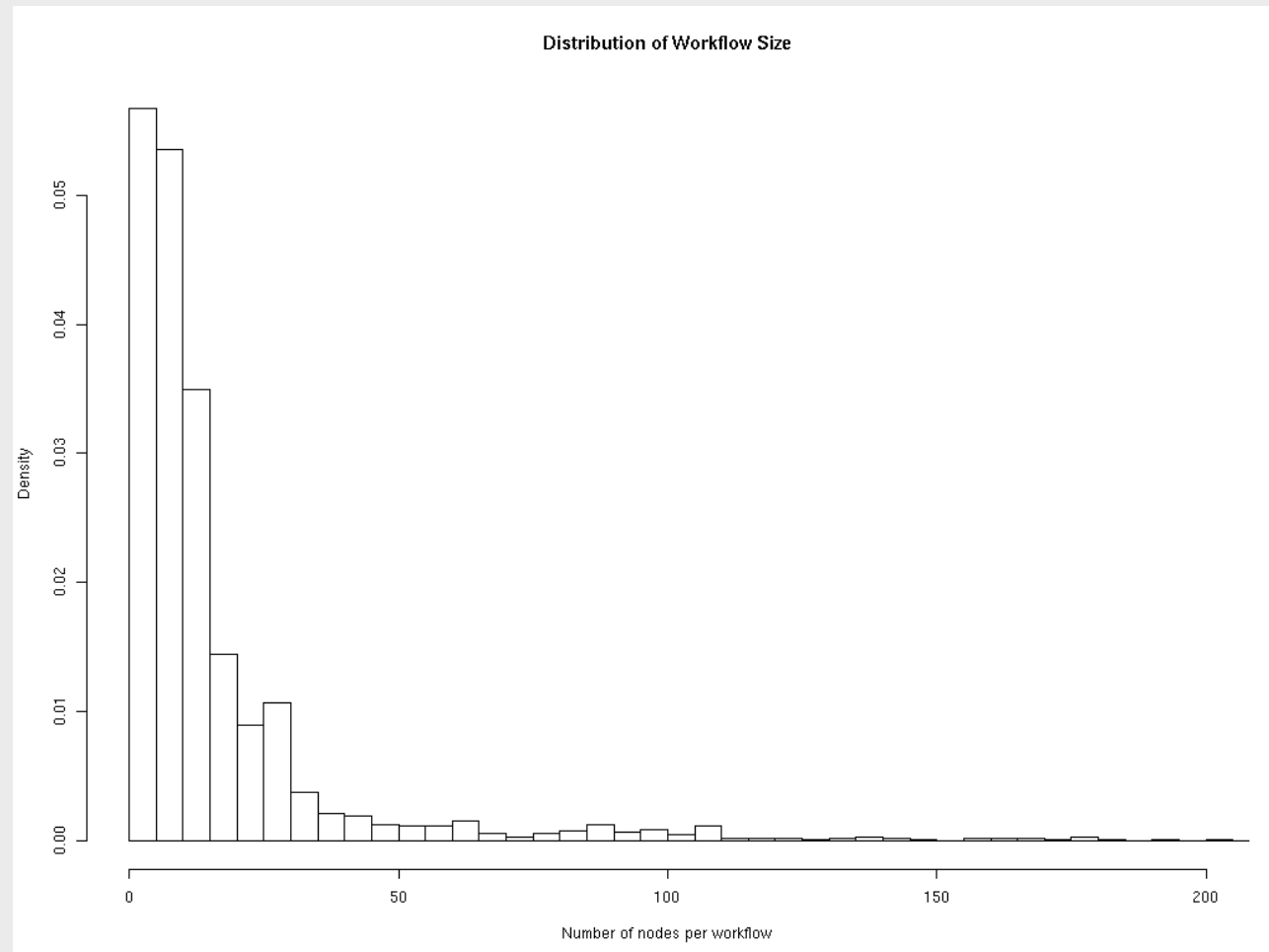
- Motivation
 - History & current status
 - Why “Workflow Mining”?
- **“Static Workflow Mining”**
 - Extraction of workflow graphs
 - Frequent sub-graph mining
 - Results
- “Dynamic Workflow Mining”
 - KNIME Profiler Plugin
 - Analysis of profiling logs
- Conclusions and Outlook



- KNIME workflows are stored in a generic XML format
- A simple parser extracts the **workflow graph**:
 - Nodes
 - Connections
- The analysis discards:
 - Node configuration / labels
 - Connection type / direction
 - Workflow layout
 - Workflow annotations
 - Personal information
 - ...

- What is frequent sub-graph mining?
 - Problem: Given a database of (labeled) graphs $\mathbf{G} = \{G_1, \dots, G_n\}$ return all (connected) sub-graphs $\mathbf{S} = \{S_1, \dots, S_k\}$ occurring at least x times in \mathbf{G} .
- Related to frequent item set mining & association rule mining
- Many applications in the life-sciences
 - Chemical structures
 - Biological networks
 - ...
- Different algorithms and implementations available
 - MoSS node in KNIME for chemical data
 - ...

- Using all workflows from central installation (~4.000 as of 2012-01-10)
- 3.899 workflows were parseable and non-empty
- 64.357 nodes
- 62.841 connections



Results - Top 10 nodes

Java Snippet



9907



Node 8

GroupBy



5275



Node 4

Joiner



4534



Node 5

Rename



4290



Node 3

Row Filter



4144



Node 11

Column Filter



3326



Node 10

File Reader



1720



Node 1

DB Open v1



1707



Node 1

Sorter



1628



Node 3

Java Edit Variable



1570



Node 2

Results - Top 11-20 nodes

Table Reader



1567



Node 4

Concatenate



1267



Node 9

Table Writer



1212



Node 5

Metanode



1136

Node 5

DB PreparedQuery v1



1051



Node 2

Reference Row Filter



893



Node 1

Double To Int



882



Node 17

DB MethodsView v1



790



Node 8

JPython Script 1:1



683



Node 3

DB Query v1



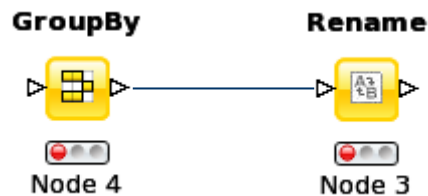
660



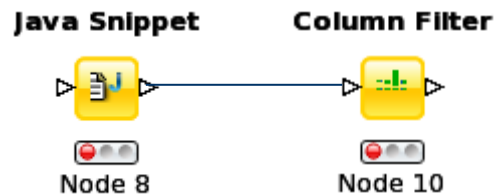
Node 6



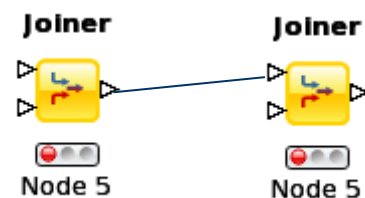
- Allow addition of multiple columns in Java Snippet ?
- Collaboration with KNIME.com



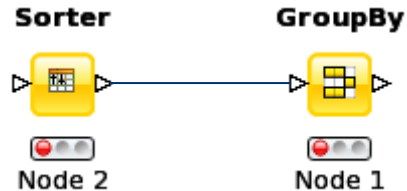
- Allow for more flexibility in column naming / typing ?



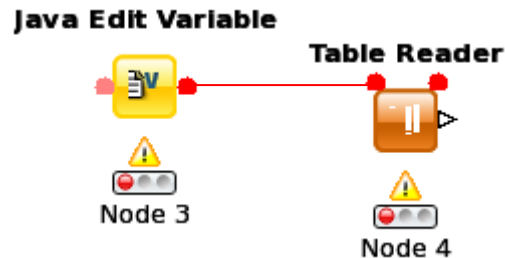
- Make users aware of “Replace columns” option



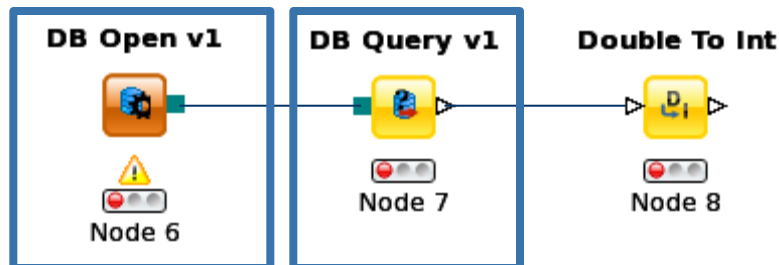
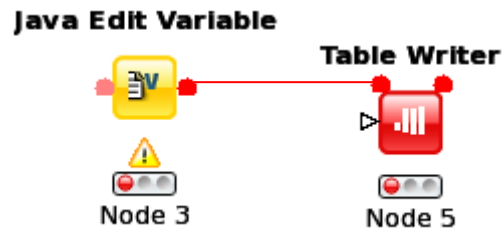
- Allow joining multiple input tables? (e.g. “Joiner (optional in)”)



➤ Issues with sorting stability until v2.3.4 (Bug 2682)



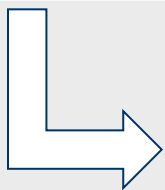
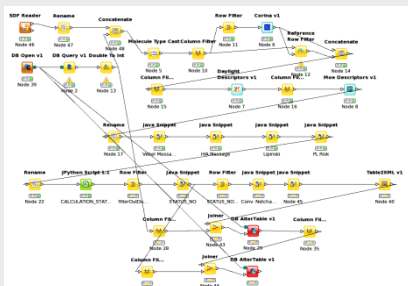
➤ Develop a concept for simple and efficient data transfer among related workflows?



➤ Improve column typing in inhouse database integration

- Motivation
 - History & current status
 - Why “Workflow Mining”?
- “Static Workflow Mining”
 - Extraction of workflow graphs
 - Frequent sub-graph mining
 - Results
- **“Dynamic Workflow Mining”**
 - KNIME Profiler Plugin
 - Analysis of profiling logs
- Conclusions and Outlook

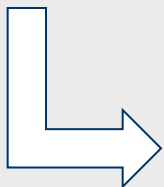
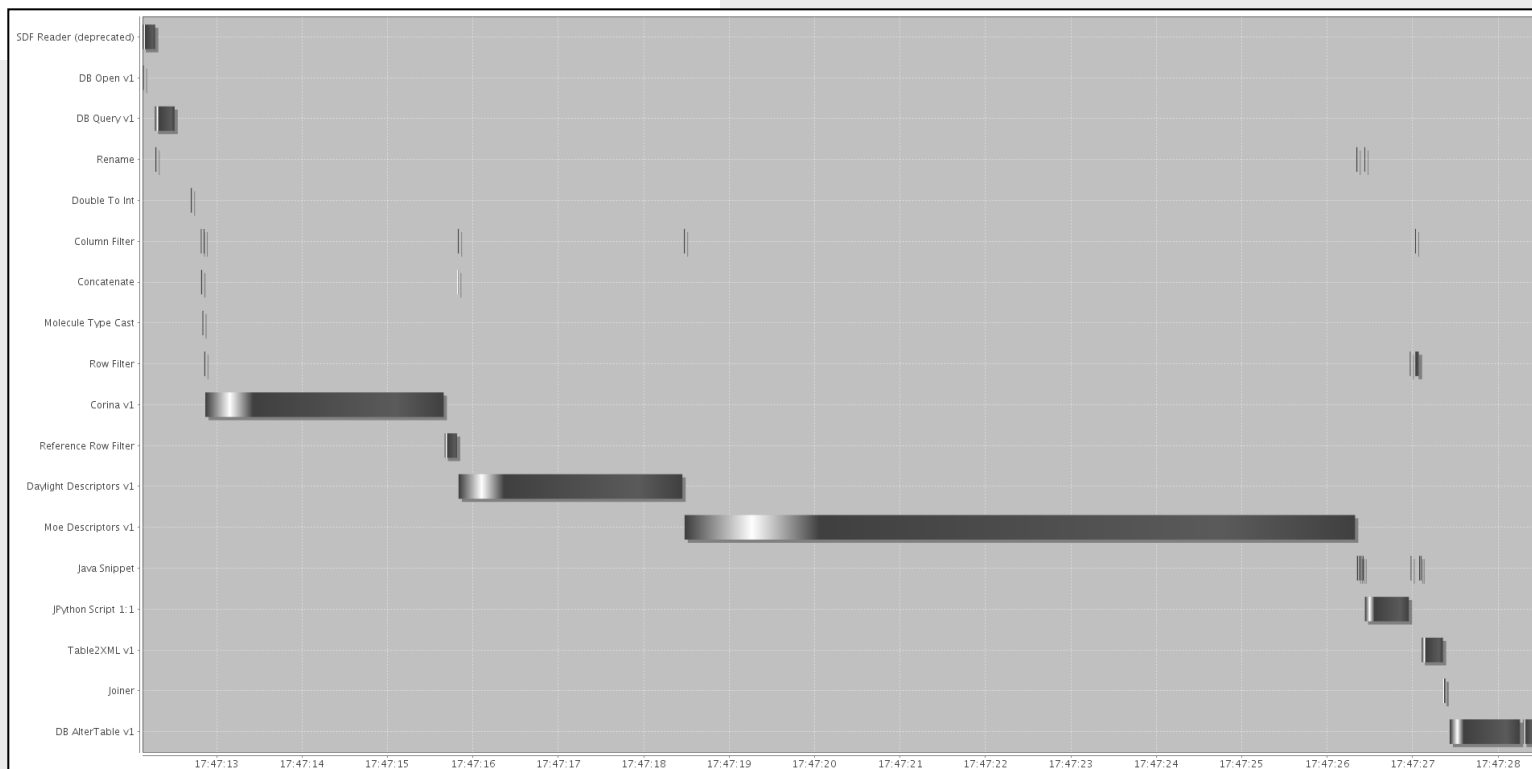
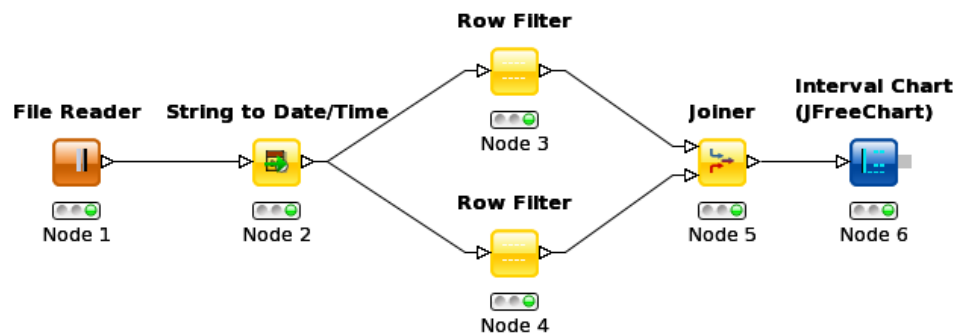
- Collection consists mainly of ad hoc-workflows
 - On average, only few executions per workflow
 - Non-interactive batch processes responsible for most workflow executions (20.000.000+ property predictions over the last 12 months)
 - Static analysis of all stored workflows might yield biased results
- No insights into bottlenecks and performance issues
 - Influence of different implementations for the same node
 - Improvements of the KNIME framework
- Limited ability to capture influence of “control flow” constructs
 - Include or exclude inactive branches in analysis?
 - How to handle loops?



Timestamp	Workflow Path	Node Name	Node Index	Custom Name	State
2011-11-12;06:45:54.662		<Metanode>	0	Node 0	EXECUTED
2011-11-12;06:45:54.670		<Metanode>	0	Node 0	CONFIGURED
2011-11-12;06:45:55.237	PROP4	DB Open v1	39	Node 39	IDLE
2011-11-12;06:45:55.240	PROP4	DB Open v1	39	Node 39	CONFIGURED
2011-11-12;06:45:55.262	PROP4	SDF Reader (deprecated)	46	Node 46	IDLE
2011-11-12;06:45:55.317	PROP4	SDF Reader (deprecated)	46	Node 46	CONFIGURED
2011-11-12;06:45:55.372	PROP4	DB Query v1	2	Node 2	IDLE
2011-11-12;06:45:55.706	PROP4	DB Query v1	2	Node 2	CONFIGURED
2011-11-12;06:45:55.751	PROP4	Rename	47	Node 47	IDLE
2011-11-12;06:45:55.757	PROP4	Rename	47	Node 47	CONFIGURED
2011-11-12;06:45:55.761	PROP4	Double To Int	13	Node 13	IDLE
2011-11-12;06:45:55.764	PROP4	Double To Int	13	Node 13	CONFIGURED
2011-11-12;06:45:55.769	PROP4	Concatenate	48	Node 48	IDLE
2011-11-12;06:45:55.774	PROP4	Concatenate	48	Node 48	CONFIGURED
2011-11-12;06:45:55.778	PROP4	Column Filter	28	Node 28	IDLE
2011-11-12;06:45:55.779	PROP4	Column Filter	28	Node 28	CONFIGURED
2011-11-12;06:45:55.790	PROP4	Molecule Type Cast	5	Node 5	IDLE
2011-11-12;06:45:55.793	PROP4	Molecule Type Cast	5	Node 5	CONFIGURED
2011-11-12;06:45:55.797	PROP4	Column Filter	10	Node 10	IDLE
2011-11-12;06:45:55.798	PROP4	Column Filter	10	Node 10	CONFIGURED
2011-11-12;06:45:55.807	PROP4	Row Filter	11	Node 11	IDLE
2011-11-12;06:45:55.807	PROP4	Row Filter	11	Node 11	CONFIGURED
2011-11-12;06:45:55.868	PROP4	Corina v1	6	Node 6	IDLE
2011-11-12;06:45:55.876	PROP4	Corina v1	6	Node 6	CONFIGURED

- Collected execution profiles for 10.000+ non-interactive batch runs
- Not used in interactive sessions, no monitoring of end-user behavior

Visualization - Gantt chart



Results - Top 10 nodes (by number of executions)

Java Snippet



161.824



Node 8

Row Filter



37.590



Node 11

Column Filter



65.282



Node 10

GroupBy



25.774



Node 4

Rename



41.138



Node 3

Reference Row Filter



24.773



Node 1

Joiner



40.672



Node 5

Double To Int



21.351



Node 8

JPython Script 1:1



38.530



Node 3

Concatenate



19.952



Node 9

Results - Top 10 nodes (by consolidated runtime)

DBBinarySimilaritySearch v1



Node 1

522.013 sec

BI Generic Model Integration v1



Node 2

419.561 sec

Moe Descriptors v1



Node 4

291.633 sec

Voisurf Descriptors



Node 5

212.338 sec

Corina v1



Node 6

138.287 sec

Blabber



Node 7

75.838 sec

Java Snippet



Node 8

66.679 sec

Daylight Descriptors v1



Node 9

58.606 sec

Table Reader



Node 10

55.274 sec

Joiner



Node 11

53.019 sec

- Motivation
 - History & current status
 - Why “Workflow Mining”?
- “Static Workflow Mining”
 - Extraction of workflow graphs
 - Frequent sub-graph mining
 - Results
- “Dynamic Workflow Mining”
 - KNIME Profiler Plugin
 - Analysis of profiling logs
- **Conclusions and Outlook**

- “Static Workflow Mining” serves as a tool to identify typical patterns in KNIME workflows
 - Supports prioritization of bug reports / feature requests
 - Identification of best practices and common pitfalls
 - Plans to provide useful Metanodes based on frequent patterns
- “Dynamic Workflow Mining” monitors workflow behavior under real-world conditions
 - Identification of bottlenecks and critical steps in workflows
 - Performance optimizations for batch-processes planned
 - Additional option to perform real-time monitoring of system health
- Large workflow collection an asset, but also a burden
 - Have to ensure compatibility of new node versions

Johannes Koppe (IS BP R&DM)

Andreas Teckentrup (IS BP R&DM)

Bernd Wiswedel (KNIME.com)

All our KNIME users for their feedback, their suggestions und their patience!