# 2nd Provenance Challenge – CESNET Team Results

*Aleš Křenek, Jiří Sitera, František Dvořák, Jiří Filipovič, Zdeněk Šustr et al.*
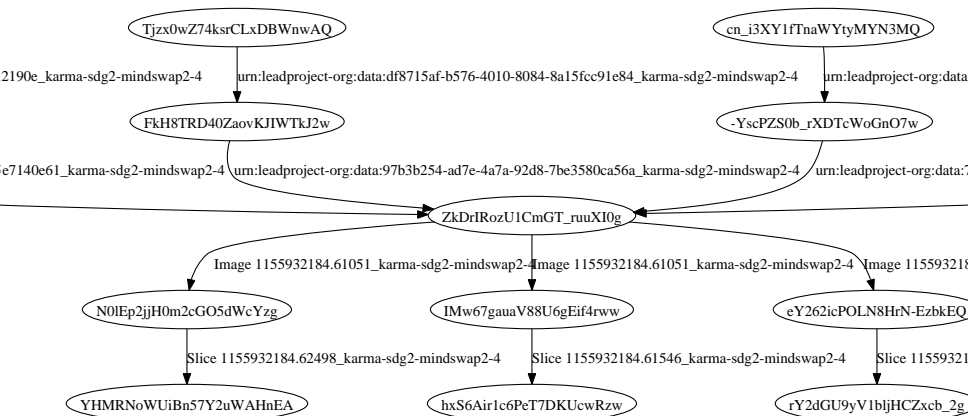
Information Society

- see our former presentations/papers for details
- queriable **archive of jobs** executed on the Grid
- job (process) is the primary entity of interest
  - all data items assigned to a particular job
- logical data view – attributes
  - namespace:name = value
- physical data store
  - individual namespace:name = value tags
  - bulk files – processed by plugins to retrive attributes
- part of EGEE gLite middleware
  - targeted at 1M jobs/day throughput

- First Challenge – explicit dependences between workflow processes
  - workflow implemented with gLite DAG
  - dependence data directly available
  - callenge query implementation depends on them
- Second Challenge – can't rely on explicit dependences
  - not present in all formats
  - definitely not available accross boundaries of split workflow parts
- detection of data dependence
  - $B$ depends on $A$ iff there is $F$ being output of $A$ and input of $B$ simultaneously
  - stored in JP explicitly $\Rightarrow$ functional query implementation
  - currently implmented as external script wrt. JP (recursive query)

- First Challenge – explicit dependences between workflow processes
  - workflow implemented with gLite DAG
  - dependence data directly available
  - callenge query implementation depends on them
- Second Challenge – **can't rely on explicit dependences**
  - not present in all formats
  - definitely not available accross boundaries of split workflow parts
- detection of data dependence
  - $B$ depends on $A$ iff there is $F$ being output of $A$ and input of $B$ simultaneously
  - stored in JP explicitely $\Rightarrow$ **functional query implementation**
  - currently implmented as external script wrt. JP (recursive query)

- convert used worklflow parts to our format
- adjust filenames (for challenge purpose only)
  - rename files to match on workflow part boundaries
  - eventually append unique suffixes to filenemes to allow multiple debugging imports of the same data
- import adjusted files
  - assign new unique job identifiers in our format
  - original ones can be presereved as process attributes
- run data-dependence detection
- run challenge queries

- five systems chosen
  - ES3, SDG, Mindswap, Karma, MyGrid
  - subjective first choice, most suitable for conversion to our format
- chosen combinations
  - CESNET-Karma-SDG, ES3-CESNET-Karma, ES3-MyGrid-SDG, MyGrid-ES3-SDG, Karma-SDG-Mindswap
- virtually **no issues specific to cross queries**
  - file naming not considered important
  - in real-world either consistent naming would be used or unique name mapping available
- all other issues related to particular system
  - running also **homogeneous** queries to demonstrate them

#1 *Find the process that led to Atlas X Graphic.*

#2 *Find the process that led to Atlas X Graphic, excluding everything prior to softmean outputting the Atlas Image.*

#3 *Find the Stage 3–5 details of the process that led to Atlas X Graphic.*

- OK in all homogeneous and heterogeneous queries
- missing warping parameters (ES)
    - not critical for these queries, used only in #1
    - affects all combinations taking ES3 part 1
- executables are named differently (#2)
    - appropriate name (matching workflow part 2) must be provided

#4 *Find all invocations of procedure align_warp that have ever occurred in the system using a twelfth order nonlinear 1365 parameter model that ran on a Monday.*

#6 *Find all images ever output from softmean where the images were align_warped using a twelfth order nonlinear 1365 parameter model.*

- depend on presence of warping parameters
  - ES3: parameters missing, query impossible
  - Mindswap: wrong format "-m -12", query possible after fixing
  - MyGrid: not implemented but present, query possible
  - SDG, Karma: query possible
- affects all combinations according to part 1

#5 *Find all Atlas Graphic images outputted from workflows where at least one of the input Anatomy Headers had an entry global maximum = 4095*
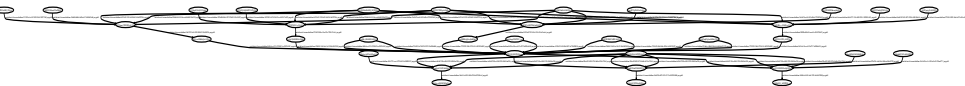
- depends on `global maximum` header entry
  - ES3, Karma, Mindswap: missing, query impossible
  - MyGrid: encoded ambiguously, not implemented
  - SDG: query possible
- affects all combinations according to part 1

#7 *A user has run the workflow twice on the same input files, in the second instance replacing each convert procedure in the final stage ...*
- done partially in First Challenge, not addressed in Second Challenge

#8 *A user has annotated some anatomy images with a key-value pair center=UChicago...*
- pure data annotation, out of scope of JP
  - but cf. MyGrid using "pseudoprocesses" producing even the inputs

#9 *A user has annotated some atlas graphics with key-value pair where the key is studyModality...*
- not implemented
- ES3, SDG, Mindswap: annotation not present, query impossible
- Karma, MyGrid: query possible

#7 *A user has run the workflow twice on the same input files, in the second instance replacing each convert procedure in the final stage . . .*
- done partially in First Challenge, not addressed in Second Challenge

#8 *A user has annotated some anatomy images with a key-value pair center=UChicago. . .*
- pure data annotation, out of scope of JP
  - but cf. MyGrid using "pseudoprocesses" producing even the inputs

#9 *A user has annotated some atlas graphics with key-value pair where the key is studyModality. . .*
- not implemented
- ES3, SDG, Mindswap: annotation not present, query impossible
- Karma, MyGrid: query possible

#7 *A user has run the workflow twice on the same input files, in the second instance replacing each convert procedure in the final stage ...*

- done partially in First Challenge, not addressed in Second Challenge

#8 *A user has annotated some anatomy images with a key-value pair center=UChicago...*

- pure data annotation, out of scope of JP
  - but cf. MyGrid using "pseudoprocesses" producing even the inputs

#9 *A user has annotated some atlas graphics with key-value pair where the key is studyModality...*

- not implemented
- ES3, SDG, Mindswap: annotation not present, query impossible
- Karma, MyGrid: query possible

Enabling Grids for E-sciencE

- processing details, file properties
- redundant data → process links
- additional processes (eg. MyGrid)



- double arcs (eg. Karma)

- "sewing" script (data dependence)
  - Challenge: standalone, invoked manually
  - production: agent, listening on JPPS feed interface for any input/output assingments
- import external formats
  - Challenge: external converter & import program, generating name=value tags
  - production: JPPS plugin, reading external format directly

- "sewing" script (data dependence)
  - Challenge: standalone, invoked manually
  - production: agent, listening on JPPS feed interface for any input/output assingments
- import external formats
  - Challenge: external converter & import program, generating name=value tags
  - production: JPPS plugin, reading external format directly

- five foreign formats successfully imported
- queries on homogeneous workflows too
- no significant differences in query results
  - clearly show format-specific issues
- virtually no issues specific to cross-queries
- local effect on possibility to perform query and query results
  - e.g. warping parameters – part 1, queries #1, 4, 6
  - cross-query results follow those on homogeneous workflows
- process vs. data provenance
  - still apparent but less serious than expected
  - some convergence seen – pseudoprocesses "generate" wf. inputs