

Analysis Preservation and Systematic Reinterpretation within the ATLAS Experiment

Kyle Cranmer, L. Heinrich
on behalf of the ATLAS collaboration

ACAT 2017



The Need for Reinterpretation

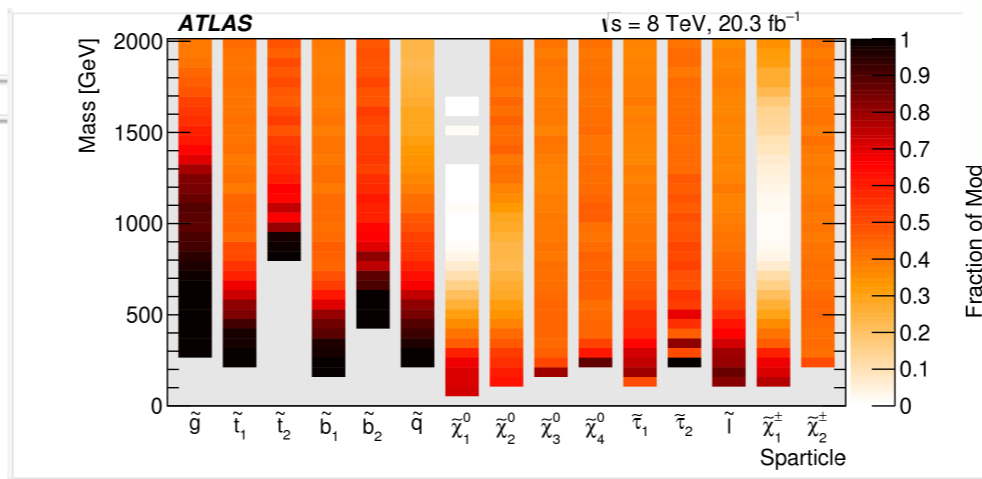
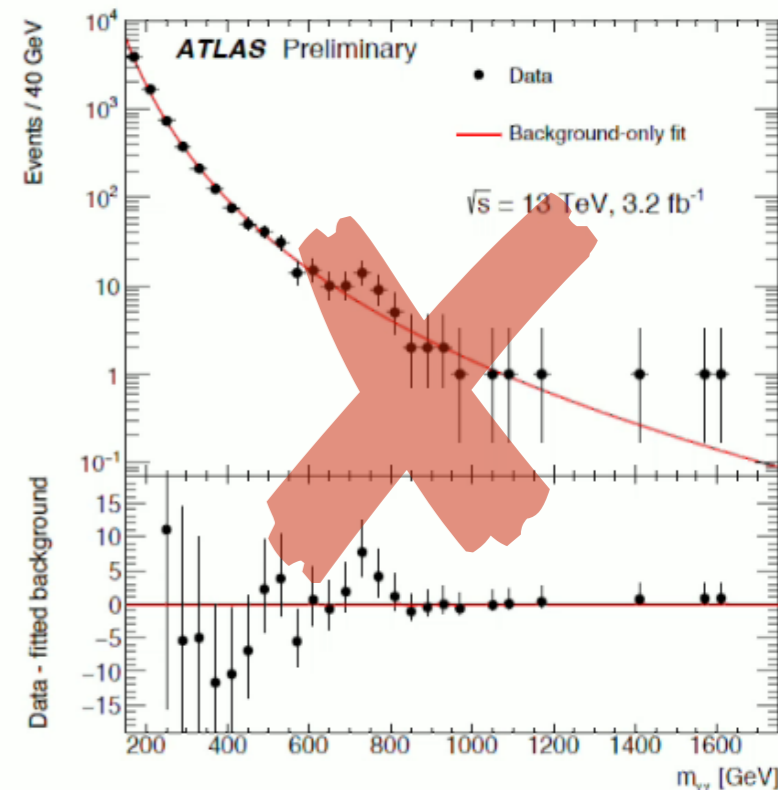
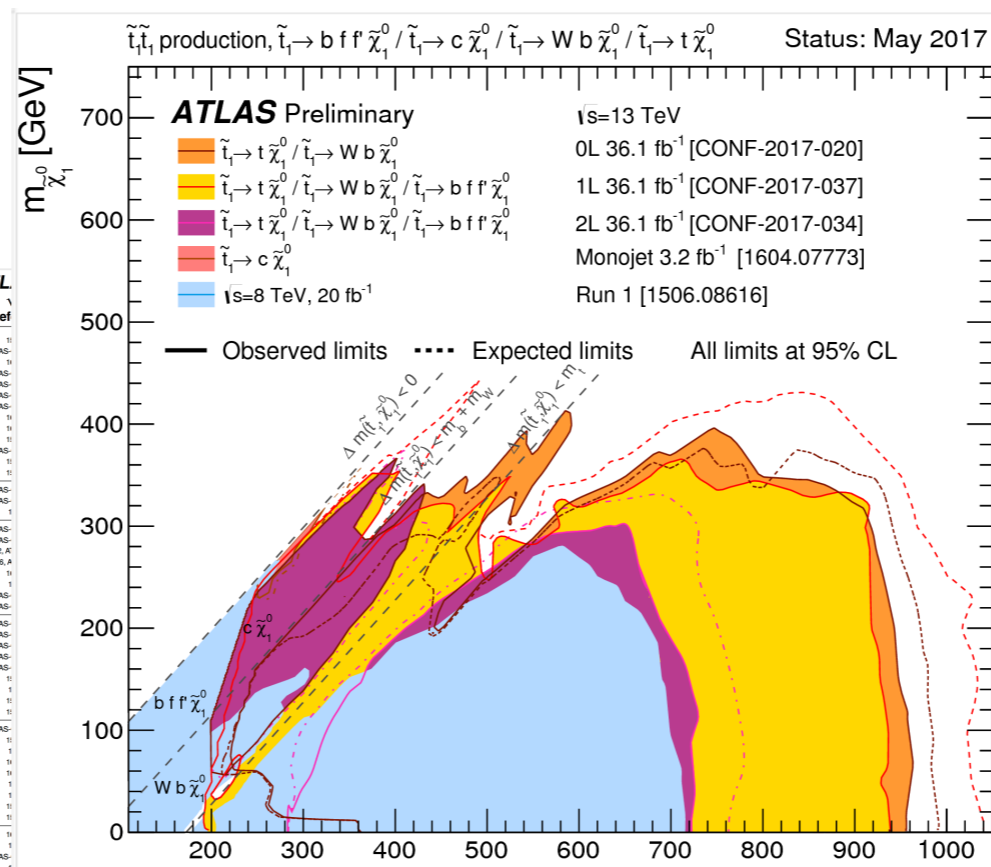
After the Higgs discovery completed the Standard Model, the search for BSM physics has become an even higher priority.

ATLAS is producing tons of results.. so far we have not found any significant excess (sorry no 750 GeV...)

ATLAS SUSY Searches* - 95% CL Lower Limits
May 2017

Model	$\epsilon, \mu, \tau, \gamma$	Jets	E_{miss} [GeV]	Mass limit	Ref
Inclusive Searches					
MSUGRA/CMSSM	0-3 $\epsilon, \mu \pm 2\tau$	2-10 jets/3b	Yes	20.3	1
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	2-6 jets	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$ (compressed)	mono-jet	1-3 jets	Yes	3.2	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	2-6 jets	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$ (with \tilde{W})	0	2-6 jets	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$ (with \tilde{Z})	3 ϵ, μ	4 jets	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$ (with \tilde{W})	0	7-11 jets	Yes	36.1	11
GMSB (\tilde{Z} NLSP)	$1.2 + 0.1 \epsilon$	0-2 jets	Yes	3.2	11
GGM (bino NLSP)	2 γ	Yes	3.2	20.3	11
GGM (higgsino-bino NLSP)	7	1 b	Yes	20.3	11
GGM (higgsino NLSP)	7	2 jets	Yes	13.3	11
GGM (higgsino NLSP)	2 ϵ, μ (Z)	2 jets	Yes	20.3	11
Gravitino LSP	0	mono-jet	Yes	20.3	11
1γ gen. production					
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	3 b	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	3 b	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	3 b	Yes	36.1	11
$\tilde{g}, \tilde{q} \rightarrow \tilde{g}^0$	0	3 b	Yes	36.1	11
3γ gen. production					
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ (SS)	1 b	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	0-2 ϵ, μ	1-2 b	Yes	4.7/13.3	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$ or \tilde{q}^0	0-2 ϵ, μ	0-2 jets/1-2 b	Yes	20.3/26.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	0	mono-jet	Yes	3.2	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ (Z)	1 b	Yes	20.3	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	3 ϵ, μ (Z)	1 b	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	1-2 ϵ, μ	4 b	Yes	36.1	11
EW direct					
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ	0	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ	0	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ	0	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	3 ϵ, μ	0	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2-3 ϵ, μ	0-2 jets	Yes	36.1	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	2 ϵ, μ, γ	0-2 b	Yes	20.3	11
$\tilde{g}, \tilde{q}, \tilde{q} \rightarrow \tilde{g}^0$	4 ϵ, μ	0	Yes	20.3	11
GGM (wino NLSP) weak prod. $\tilde{q}^0 \rightarrow \tilde{q}^0$	1 $\epsilon, \mu + \gamma$	0	Yes	20.3	11
GGM (bino NLSP) weak prod. $\tilde{q}^0 \rightarrow \tilde{q}^0$	2 γ	0	Yes	20.3	11
Long-lived particles					
Direct \tilde{q}^0 prod., long-lived \tilde{q}^0	Disapp. trk	1 jet	Yes	36.1	11
Direct \tilde{q}^0 prod., long-lived \tilde{q}^0	DE/dx trk	Yes	18.4	11	
Stable, stopped \tilde{q}^0 hadron	0	1-5 jets	Yes	27.9	11
Stable \tilde{q}^0 hadron	trk	-	-	3.2	11
Metastable \tilde{q}^0 hadron	DE/dx trk	-	-	18.1	11
GMSB, stable \tilde{q}^0 hadron	1 μ, γ	-	Yes	15.1	11
GMSB, \tilde{q}^0 long-lived \tilde{q}^0	2 γ	-	Yes	20.3	11
GGM $\tilde{q}^0 \rightarrow \tilde{q}^0$	disapp. trk	-	-	20.3	11
GGM $\tilde{q}^0 \rightarrow \tilde{q}^0$	disapp. trk	-	-	20.3	11
RPV					
LFV $\tilde{q}^0 \rightarrow \tilde{q}^0 + \tilde{q}^0 \rightarrow \tilde{q}^0$	$q_1 q_2 q_3$	-	-	3.2	11
Bilinear RPV CMSSM	2 ϵ, μ (SS)	0-3 b	Yes	20.3	11
$\tilde{q}^0 \rightarrow \tilde{q}^0 + \tilde{q}^0$	4 ϵ, μ	-	Yes	13.3	11
$\tilde{q}^0 \rightarrow \tilde{q}^0 + \tilde{q}^0$	3 $\epsilon, \mu + \tau$	-	Yes	20.3	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	0	4-5 large-R jets	Yes	14.8	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	0	4-5 large-R jets	Yes	14.8	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	1 ϵ, μ	8-10 jets/0-4 b	Yes	36.1	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	1 ϵ, μ	8-10 jets/0-4 b	Yes	36.1	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	0	2 jets + 2 b	Yes	15.4	11
$\tilde{q}^0 \rightarrow \tilde{q}^0$	2 ϵ, μ	2 b	Yes	36.1	11
Other					
Scalar charm, $\tilde{c} \rightarrow \tilde{c}^0$	0	2 c	Yes	20.3	11

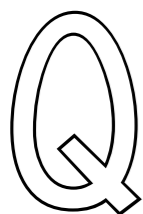
*Only a selection of the available mass limits on new states or phenomena is shown. Many of the limits are based on simplified models, c.f. refs. for the assumptions made.



The Need for Reinterpretation

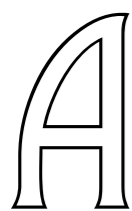
Where is the New Physics?

- hide in unexpected places, complex final states, low-rate / low-acceptance scenarios (e.g. compressed models, models spreading across many topologies)
- not be reachable at all at the LHC



how do we exploit the LHC data to maximize our understanding of the model landscape?

there are **many more** candidate models than we have graduate students to design dedicated analyses for each new model — let's make the most of the analyses that we do have. **Many of them are sensitive to a whole range of models.**



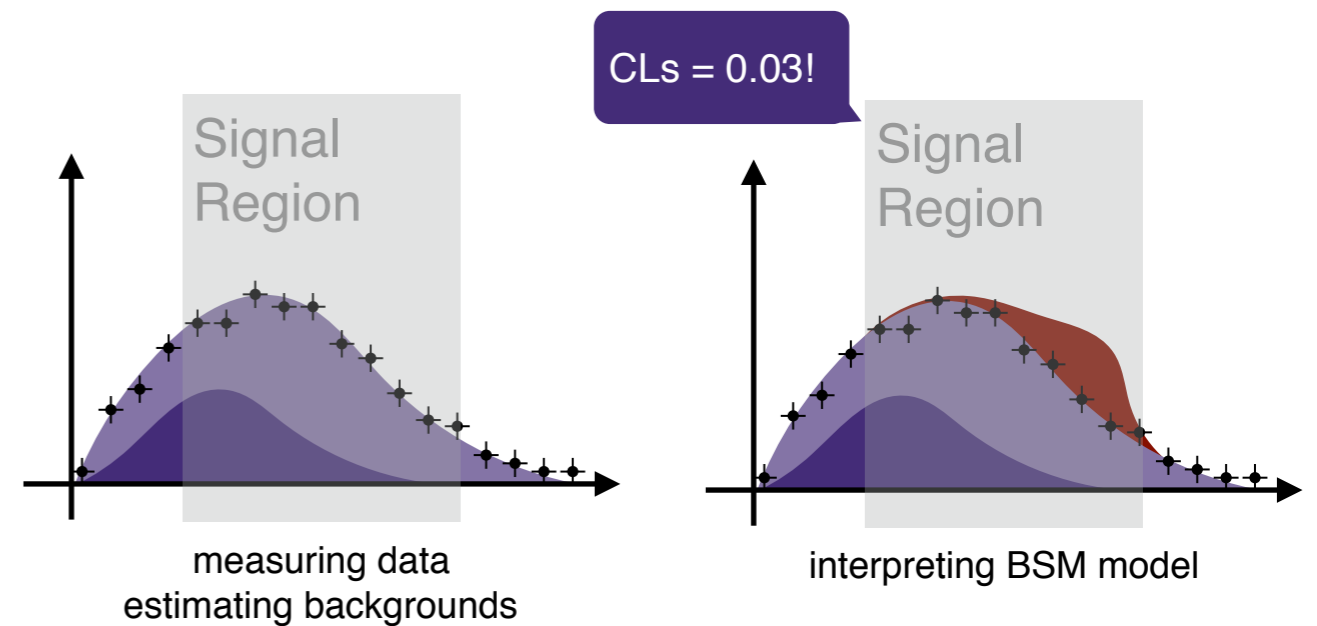
leverage modern analysis preservation and reusability techniques to *re-interpret existing analyses*



The Need for Reinterpretation

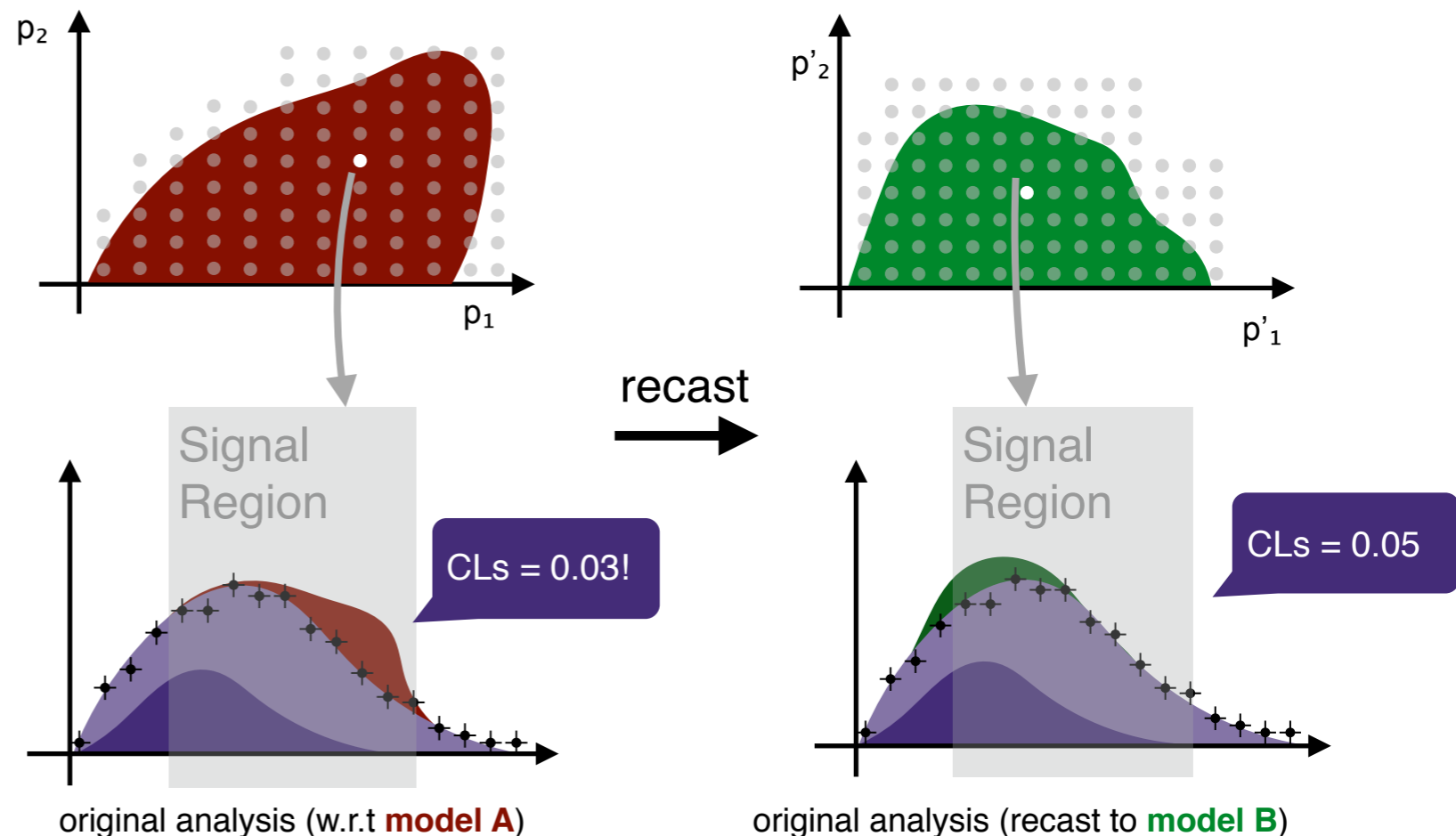
Most of the work goes into: **taking data, designing, validating** the analysis strategy, **understanding Standard Model backgrounds**.

Model interpretation come at the end, and are technically the **easiest part**: analysis pipeline is **fixed** after unblinding, MC dataset sizes small. Analysis teams routinely check hundreds of parameter points (of their favorite model).



Reinterpreting / Recasting:

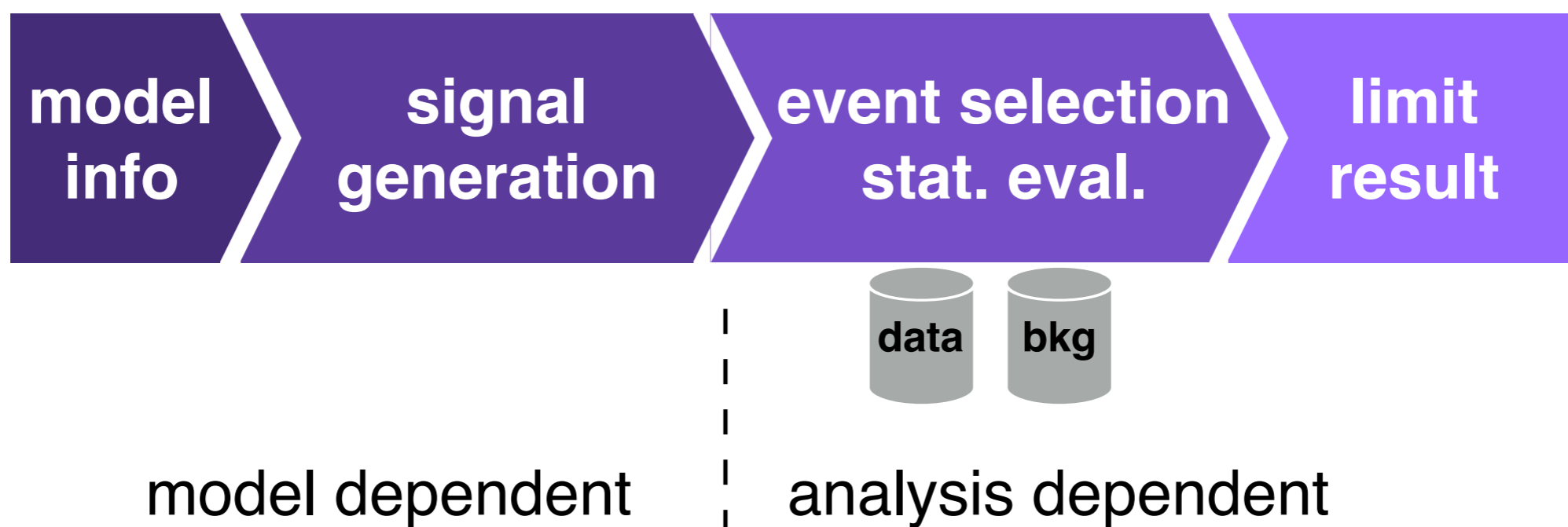
1. predicting BSM contributions of new model under a given analysis (**same** event selection)
2. statistical analysis of new signal with respect to **same** data and background estimates as original analysis to derive **new limits**



The Recipe for Reinterpretation

Reinterpretation follows a straight-forward recipe with three ingredients:

1. Ability to generate new signal model (incl. access to DetSim, Reco)¹
2. Access to the analysis / event selection logic
3. Access to data and background distributions (incl. systematic variations) for statistical analysis

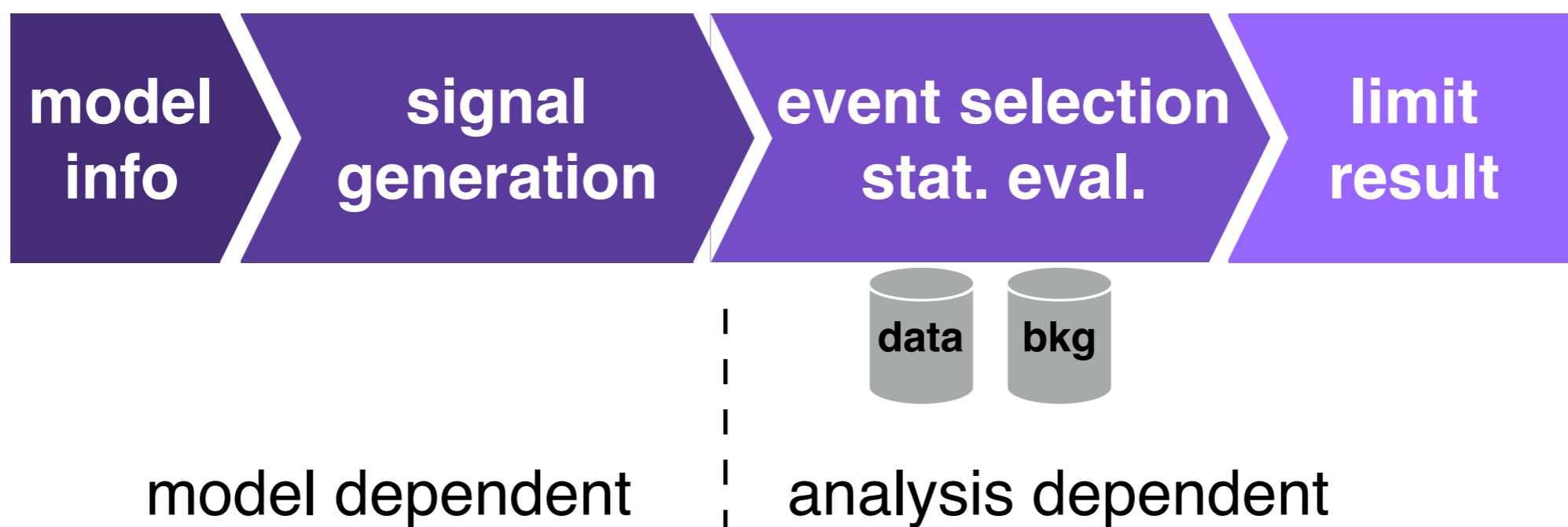


The Recipe for Reinterpretation

Reinterpretation follows a straight-forward recipe with three ingredients:

1. Ability to generate new signal model (incl. access to DetSim, Reco)¹
2. Access to the analysis / event selection logic
3. Access to data and background distributions (incl. systematic variations) for statistical analysis

- true implementation of recipe requires access to collaboration-internal data/software
- there is an eco-system built by the pheno community to approximately implement the recipe CheckMate, ATOM, SModelS, SUSY-AI, etc..

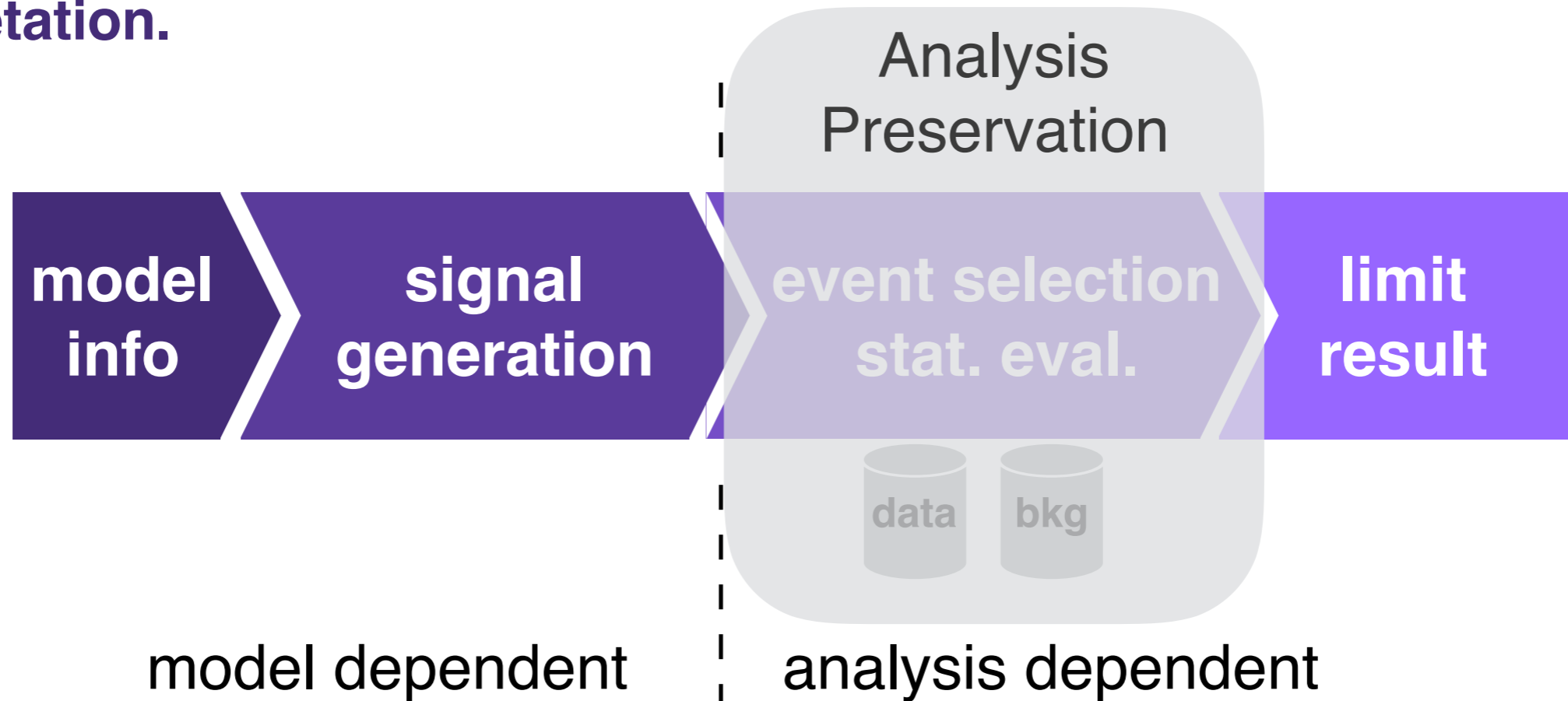


Reinterpreting in ATLAS

ATLAS has well-managed and standardized process to generate new signals — almost anyone can do it / request new samples — not the bottleneck.

The most challenging aspect for ATLAS has been the preservation of the downstream analysis code

Solving analysis preservation / reusability enables systematic reinterpretation.



Challenges for analysis preservation

- need to preserve such that it can be **re-run on new input**
- real ATLAS analyses are complex. Not a single file in a common framework (like e.g. Rivet, CheckMate, LHADA). **There's a reason have our own computing model.**
 - code is very diverse. many frameworks, scripts, etc..
- distributed teams, code, data: **one person rarely is able to run the entire analysis** pipeline — some develop event selection, some background estimates, some statistical analysis

To preserve analyses, we needed to respect the tools, workflows people use. instead of forcing a re-implementation, develop toolchain to capture what they are already doing.

1. capture software (*including all dependencies*) needed to run individual parts of an analysis (e.g. event selection) in a future-proof way.
2. capture logic how the many pieces of the analysis fit into an *analysis workflow* that can be re-executed on a new signal



Analysis Preservation in ATLAS

comprehensive software capture was intractable until recently (VMs??). Now progress in IT industry has **made it feasible** — **Linux Containers**. Technology with wide industry support — will be here for foreseeable future.



revolutionized software distribution & archival — “app store for generic software”. Many additional tools that help deploy / run Linux Containers in “the cloud” (Google, Amazon, Microsoft, etc...).

Containers are now becoming a major topic in LHC collaborations. Simplifies a lot of our computing in many ways.

technology stack enabling realistic analysis preservation has become available recently



ceph



CernVM
File system



kubernetes



docker



openstack™



Containerizing common ATLAS software:

centralized effort by software infrastructure team to build base images on which standard ATLAS software (e.g. offline reconstruction, analysis frameworks) run. (see poster by M .Vogel)

Choice between “lean” and “fat” containers

- provides software preservation for common ATLAS workloads (reconstruction, simulation, etc)
- provides easy-to-use base images for end-users



Analysis Preservation in ATLAS

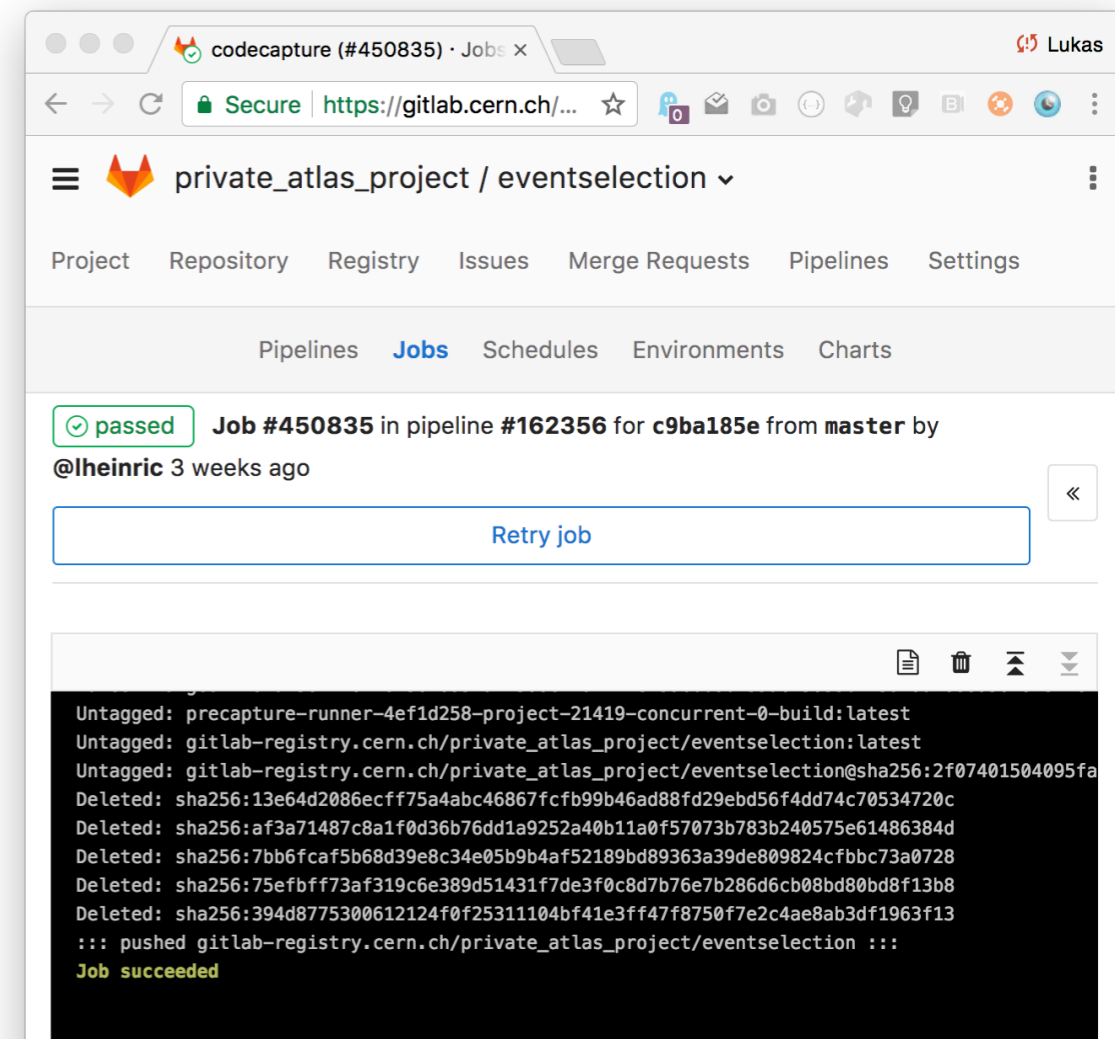
Enabling User-Level Analysis Software Preservation via Continuous Integration

Crucial to have streamlined system to capture user analysis code into containers — not everyone will become container expert.

Increasing usage of CERN GitLab installation for managing user code allows us to use built-in continuous integration infrastructure to capture code **at development time. Avoids asking people for code after publication.**

```
.gitlab-ci.yml 219 Bytes
1 codecapture:
2   tags: [code-capture]
3   variables:
4     BUILD_BASE_IMAGE: lukasheinrich/recast_cvmf:
5     image: gitlab-registry.cern.ch/codecapture/co
6   script:
7     - /codecapture_utils/steering.sh

docker.build.sh 98 Bytes
1 echo "building analysis code..."
2
3 rcsetup Base,2.4.14
4 rc find_packages
5 rc compile
6
7 echo "done.. " 11
```



Analysis Preservation in ATLAS

Declarative Parametrized Workflow Description

developed declarative workflow language *yadage* based on ubiquitous industry standards (JSON) to describe logic between separate analysis stages.

allows arbitrary, runtime-dependent, directed acyclic graphs of containerized analysis workflows

models original analysis workflow instead of forcing common interface.

Can be developed, validated *during analysis development*, stored as data fragment in repo

parameters to bind
new inputs
(*samples + xsec*)

```
stages:
- name: selection
  dependencies: ['init']
  scheduler:
    scheduler_type: singlestep-stage
    parameters:
      cxaodfile: {stages: init, output: cxaodfile, unwrap: true}
      did:       {stages: init, output: did, unwrap: true}
      xsec_pb:   {stages: init, output: xsec_pb, unwrap: true}
      k_factor: {stages: init, output: k_factor, unwrap: true}
      filt_eff: {stages: init, output: filt_eff, unwrap: true}
      nametag: 'recast'
      mctype:   {stages: init, output: mctype, unwrap: true}
      outputdir: '{workdir}/output'
      step: {$ref: 'selscript.yml#'}
- name: fit
  dependencies: ['selection']
  scheduler:
    scheduler_type: singlestep-stage
    parameters:
      fitinputdir: '{workdir}/fitinputs'
      signalfile: {stages: selection, output: out, unwrap: true}
      zerolepbkg: 'root://eosuser.cern.ch///eos/project/r/recast/at
      onelepbkg:  'root://eosuser.cern.ch///eos/project/r/recast/at
      twolepbkg: 'root://eosuser.cern.ch///eos/project/r/recast/at
      nametag: 'recast'
      limitfile: '{workdir}/limits.txt'
      plotdir: '{workdir}/plots'
      step: {$ref: 'fitscript.yml#'}
```

```
selscript.yml 974 Bytes
1 process:
2   process_type: 'interpolated-script-cmd'
3   interpreter: bash
4   script: |
5     source ~/.bashrc
6     source ./rcSetup.sh
7     /recast_auth/getkrb.sh
8
9     echo {cxaodfile} > /tmp/filelist
10    cat FrameworkSub/data/XSections_13TeV.txt|grep -v '^{did}'> tmp.txt
11    mv tmp.txt FrameworkSub/data/XSections_13TeV.txt
12    echo 'adding this line to XSections_13TeV.txt'
13    echo '{did} {xsec_pb} {k_factor} {filt_eff} {nametag} {mctype}_RECAST'
14    echo '{did} {xsec_pb} {k_factor} {filt_eff} {nametag} {mctype}_RECAST' >> F
15    hsg5frameworkReadCxADD_monoVH {outputdir} data/FrameworkExe_monoVH/framework
16 publisher:
17   publisher_type: interpolated-pub
18   publish:
19     out: '{outputdir}/hist-testrun.root'
20 environment:
21   environment_type: 'docker-encapsulated'
22   image: gitlab-registry.cern.ch/priecck/monohbb16_preservation/eventselection
23 resources:
24   - CVMFS
25   - GRIDProxy
```

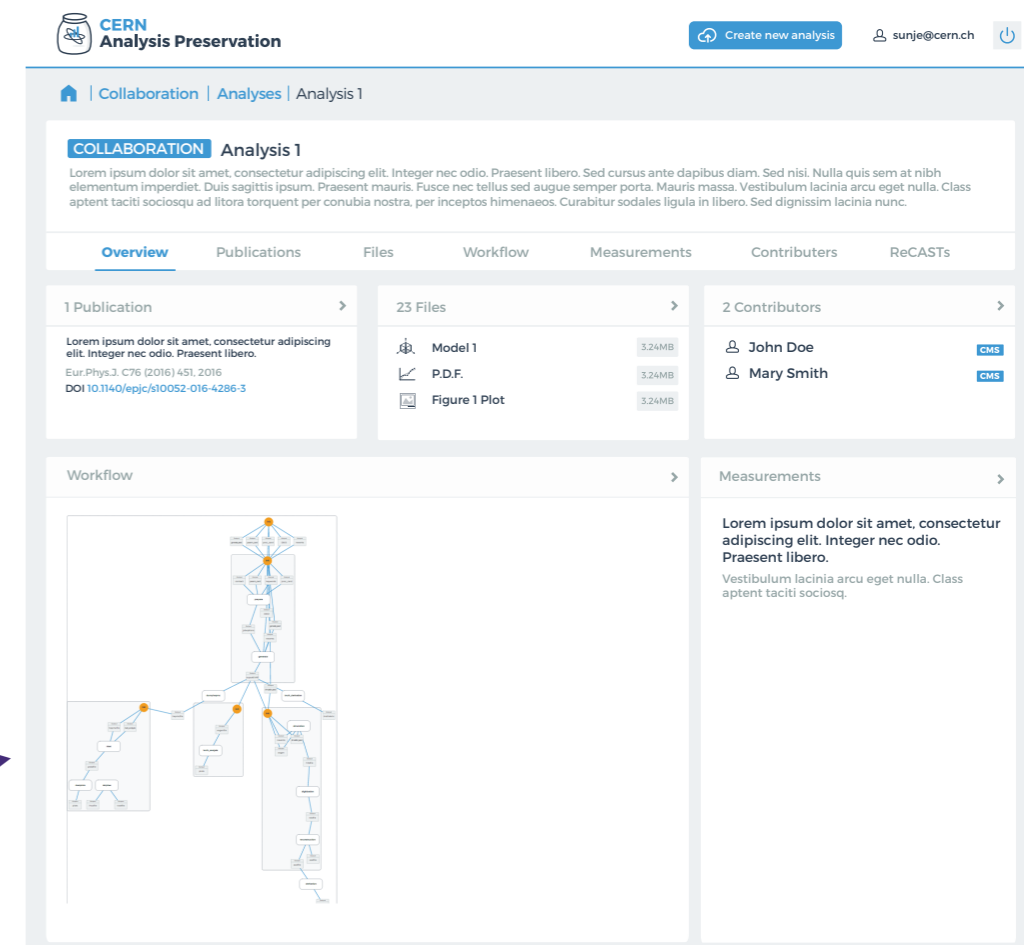
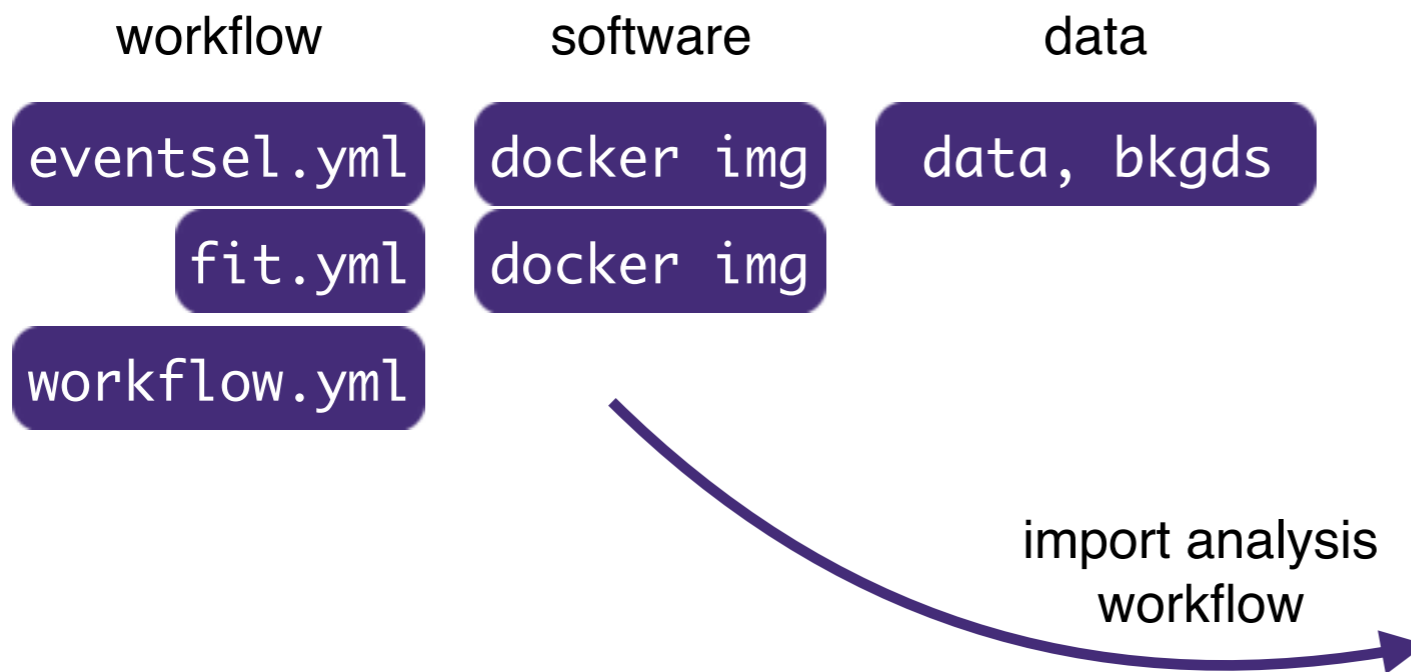
referencing
archived assets
(data, bkg histos)



Analysis Preservation in ATLAS

Close collaboration with **CERN Analysis Preservation (CAP) Portal**. Invenio-based portal that can natively ingest:

- code and software environments (clone repos, import docker images)
- native support for workflow specification (yadage workflows adhere to JSON schema spec)
- ingestion of data assets (data + background histograms / ntuples, etc..)



to recap..

analysis preservation during development

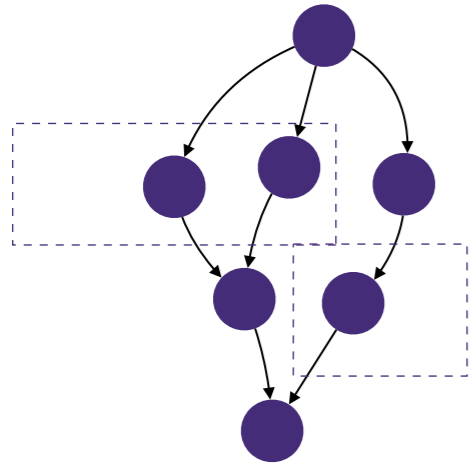
**capture analysis software using linux containers
through continuous integration**

**standards-based workflow definition to build
workflow graph of analysis stages**

deep integration with CERN Analysis Preservation



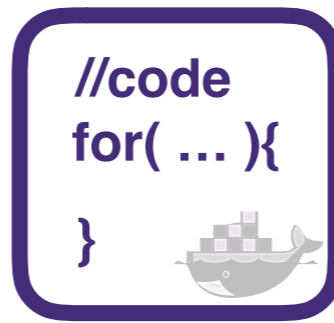
workflows



data



code + env



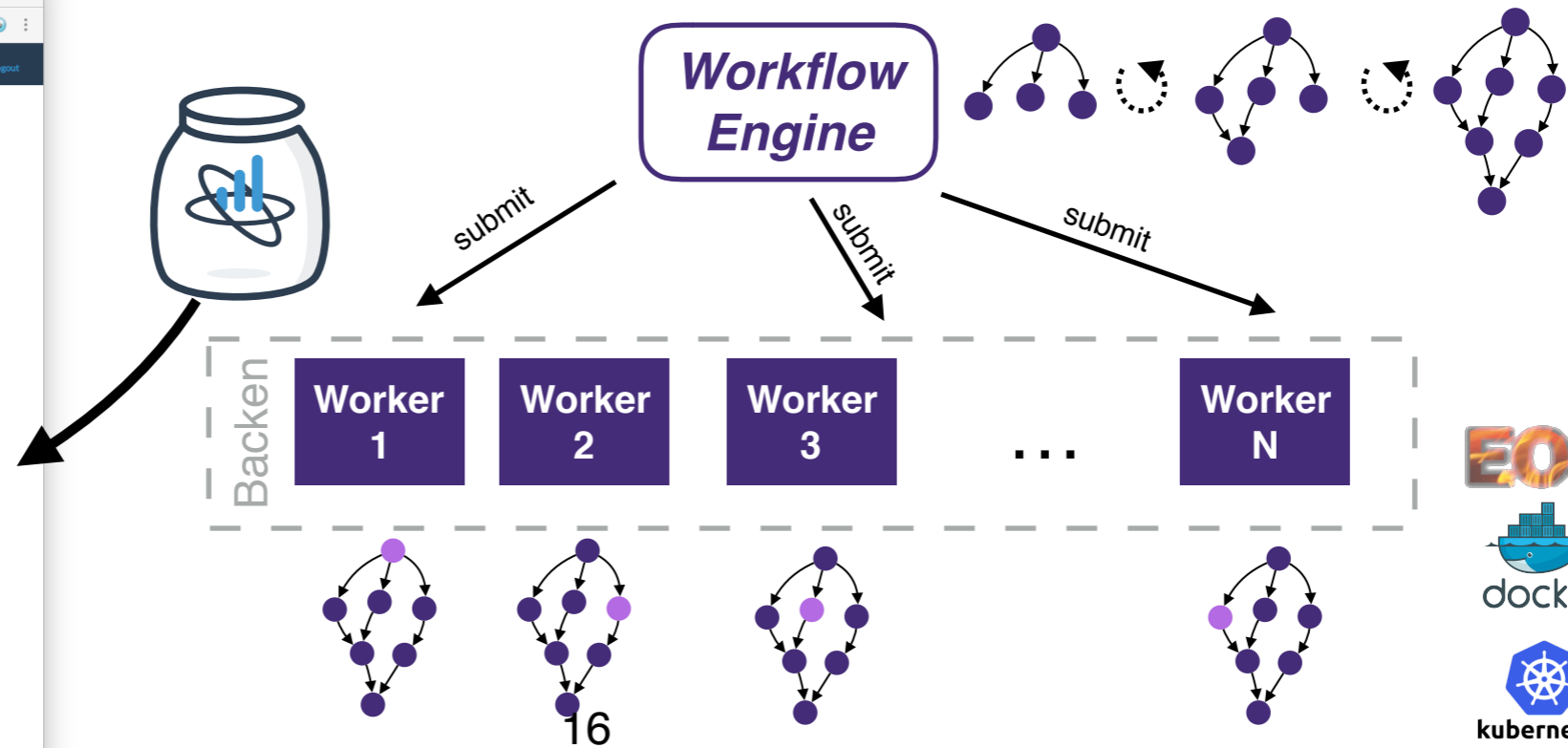
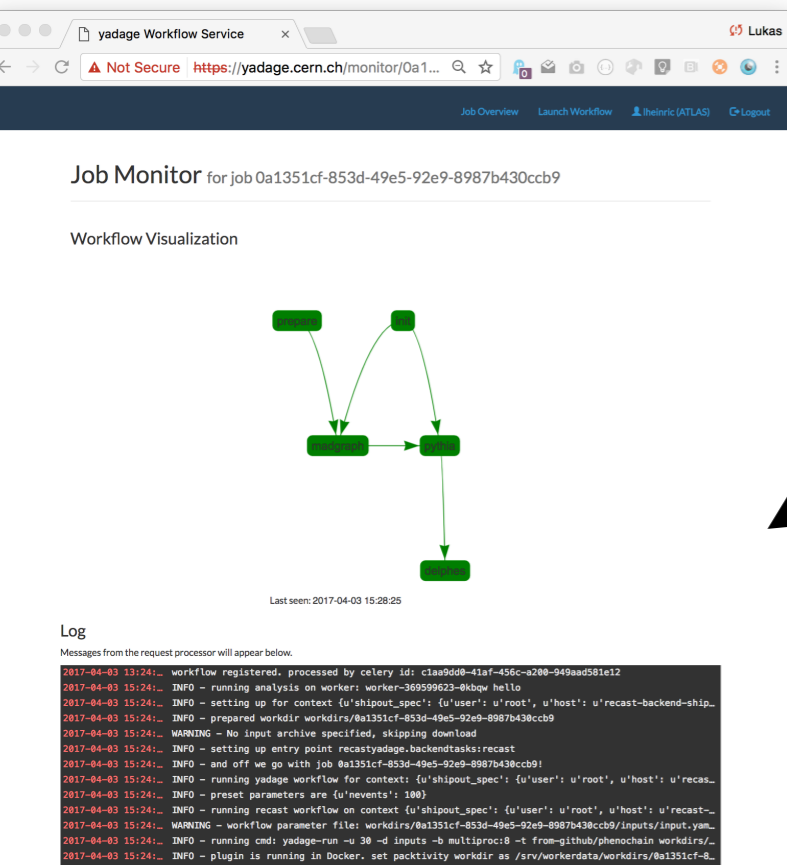
Once archived, we can build new client applications by re-using analyses



REANA

REANA — generic Workflows-as-a-Service platform

- main unit of work is a containerized workflow as stored in CAP
- “workflow engines” control what container jobs to submit
 - yadage support at launch
 - exploring Common Workflow Language
- deployed natively on Kubernetes. Deployable at AWS, GCE, CERN OpenStack.
- Multi-experiment support (tested w/ ATLAS, LHCb)
- Joint Effort by CERN IT and SIS, DASPOS, DIANA-HEP, funded by NSF, MSDSE

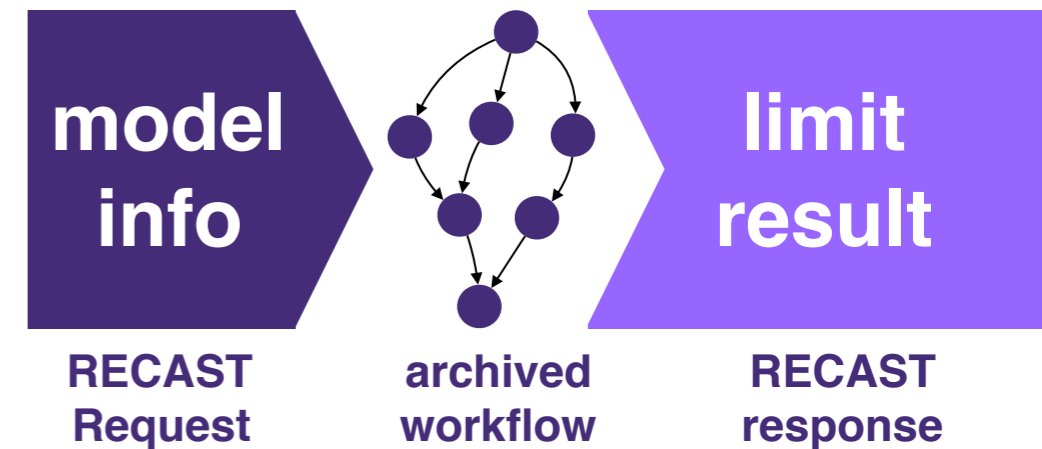


RECAST

RECAST — a semantic layer on top of REANA for reinterpretations

Reinterpretations provide a fundamental interface that can be satisfied by a (combination of) containerized workflows.

Idea:



provide web-based interface to request and process reinterpretations powered by preserved analyses in CAP

- produce reinterpretations of same fidelity as original result (not just approximations)
- allow users to request new reinterpretations, provide parameter points, model data (SLHA files, etc.)
- control interface to manage and view reinterpretation, fulfillment progress, submission API
- Integration with publishing APIs (HepData, Inspire, etc.)



RECAST

RECAST – Infrastructure

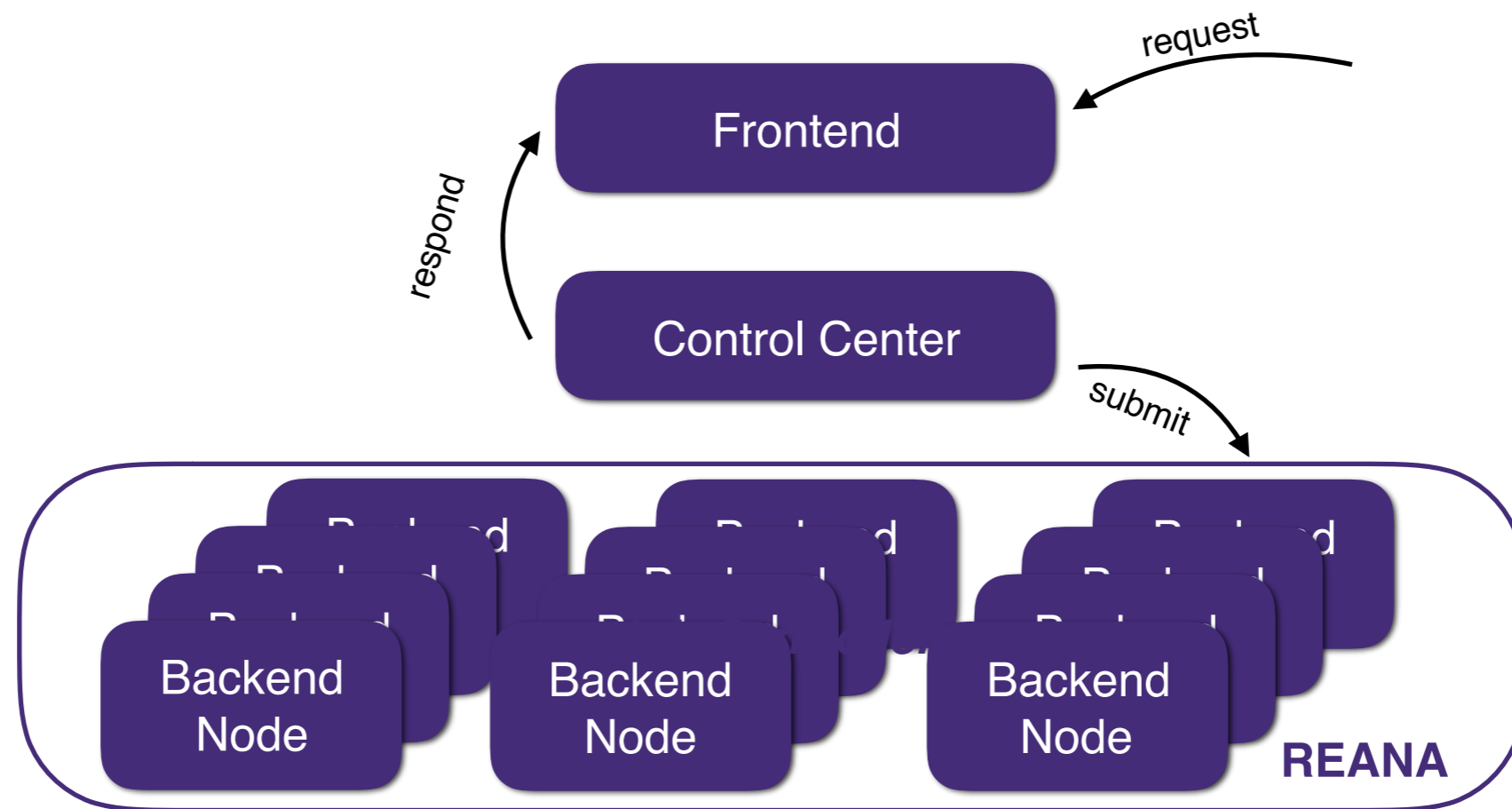
Frontend: register requests and analyses, manage request data, show reinterpretation results

Control Center: web interface to view and process requests, manage implementation catalogue.

Backend Cluster: REANA workflow service.

Deployed on CERN infrastructure

Eco-system of tools and libraries: REST API, python bindings, command line tools.



RECAST

Importing a new analysis:

```
$> recast import_analysis from_web cds/1525880
```

import via API import from ArXiv, Inspire, CDS (for CONF notes)

Creating new Parameter Scans:

```
$> recast createscan scan.yml
```

provide list of parameter points
and add concrete request data
as zip file for each point
(format specified by unique identifier)

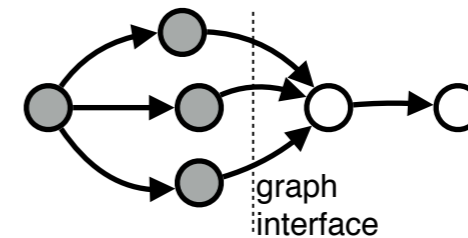
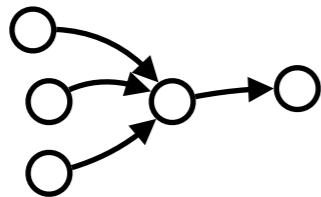
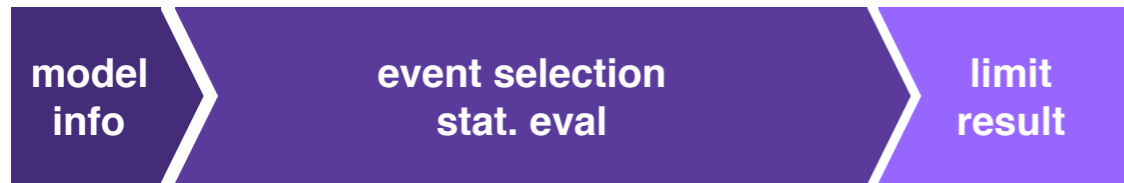
```
type: recast_scan
title: 'Validation of ATLAS-CONF-2013-024 in Neutralino/Stop Mass Plane'
pubkey: 'cds/1525880'
request_format: 'standard_format'
description: >
  this reinterpretation re-validates the original grid in of ATLAS-CONF-2013-024
reason: >
  The grid is useful to validate / check third-party
  (such as CheckMate) implementations of the ATLAS-CONF-2013-024 as it
additional_information: >
  the grid consists of 36 points in the stop / neutralino mass plane.
  The input parameters are given in the standard format consisting of
  model parameters and number of events.parameters: [mStop, mNeutralino]
parameters: [mStop, mNeutralino]
points:
- coordinates: [200.0, 0.0]
  data: data/200.0_0.0.zip
- coordinates: [300.0, 0.0]
  data: data/300.0_0.0.zip
- coordinates: [300.0, 100.0]
  data: data/300.0_100.0.zip
  ...
```



RECAST

Processing Requests using archived workflows:

if model info (request) data already in form that analysis workflow expects, can just process in a single-pass

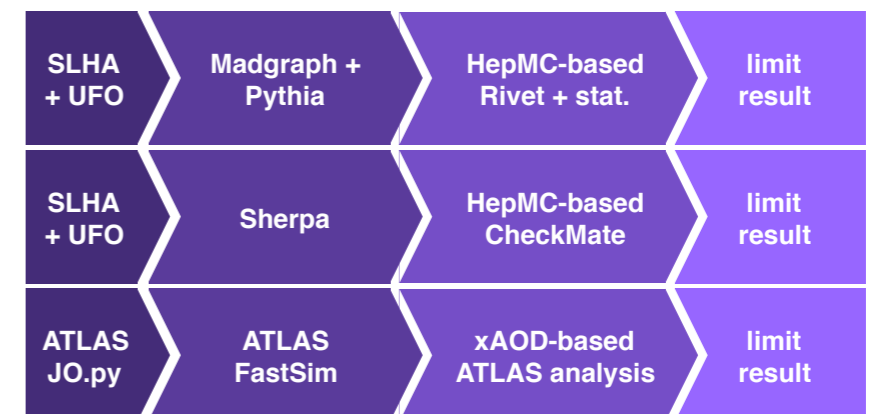
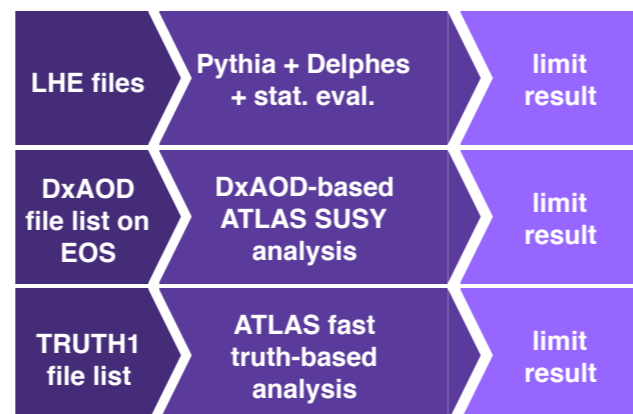


Otherwise, we can use *graph interfaces* between workflows to build more complex workflows by **mixing and matching** separately preserved workflows.

- possible, because yadage workflows are machine readable. can programmatically create new workflows.

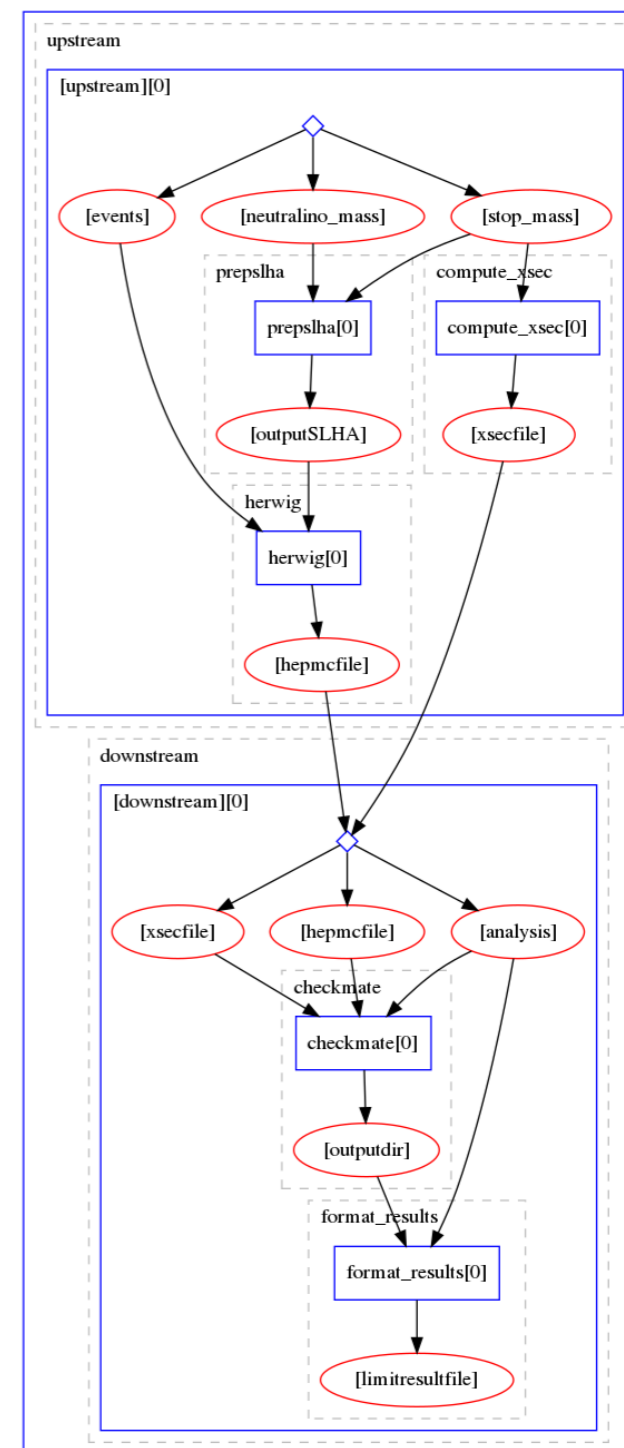
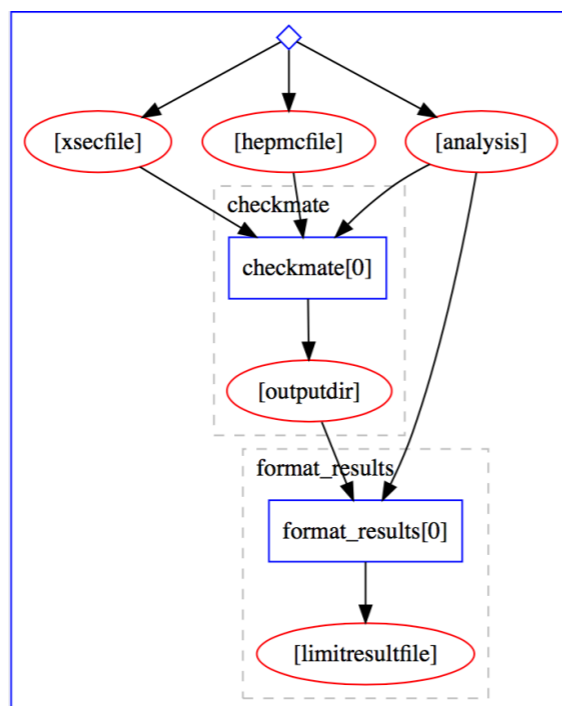
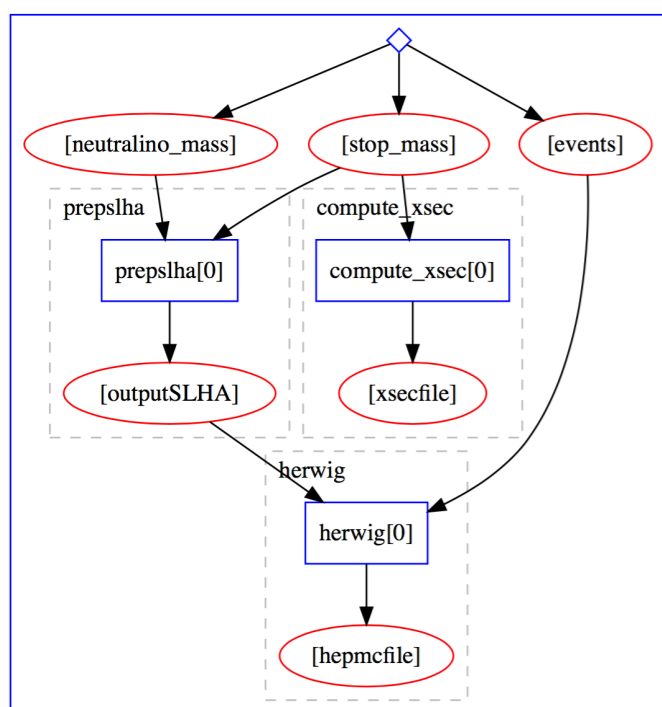
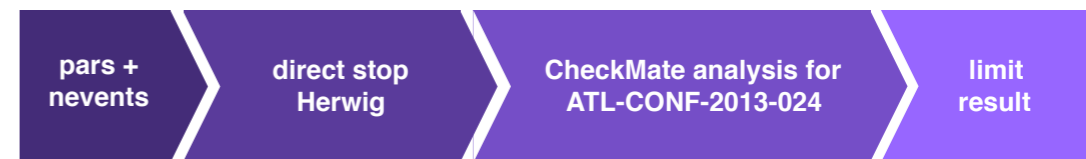
Allows us to build up a **comprehensive catalogue** of reinterpretation workflows quickly.

- match entire CheckMate/Rivet catalogues with generic HepMC producing workflows.
- only require analysis teams to capture their immediate workflow based on derived xAOD ntuples, upstream **tuple production workflow can be handled separately.**



RECAST

Example: Herwig upstream with CheckMate downstream



Herwig workflow

CheckMate workflow

computed combined workflow

upstream and downstream workflow can be developed and maintained independently, on mutually exclusive software stack. Combination works on the semantic workflow level / between files on disk.



RECAST

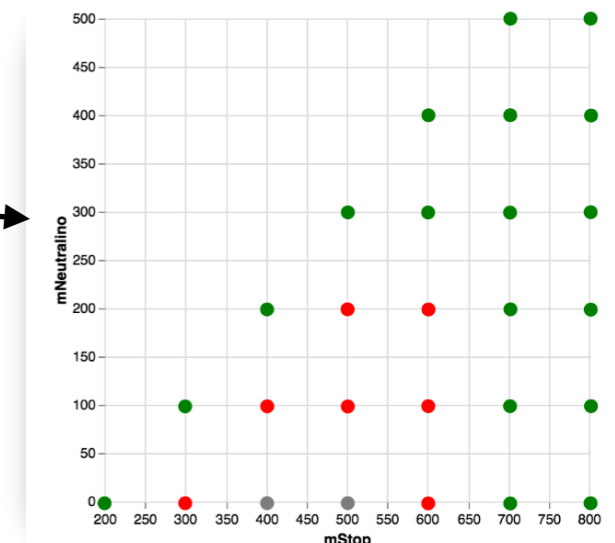
Processing Requests from the Control Center

web interface showing request details and result visualization. Offers all possible workflow configurations given request format and analysis details.

Request Details
© 2017-08-13
Analysis 4
Search for direct production of the top squark in the all-hadronic ttbar + etmiss final state in 21 fb-1 of p-p collisions at sqrt(s)=8 TeV with the ATLAS detector
Reason:
The grid is useful to validate / check third-party (such as CheckMate) implementations of the ATLAS-CONF-2013-024 as it
Additional Info:
the grid consists of 36 points in the stop / neutralino mass plane. The input parameters are given in the standard format consisting of model parameters and number of events.

Requested Parameter Points
Parameter Point 1
1 Basic Requests mStop: 200.0 mNeutralino: 0.0 obs: 1.00 exp: 1.00

toggle scan scatter plot between Cls heat map and 0.05 exclusion threshold



Process Show Processings Results Upload

- requestwflow-checkmate (None)
- stops_herwig_nllfast-atlas_analysis (standard_format)
- stops_herwig_nllfast-checkmate (standard_format)
- madgraph_pythia-checkmate (nevents_run_pars)
- requestwflow-atlas_analysis (None)
- madgraph_pythia-atlas_analysis (nevents_run_pars)

possible processing workflows for this request format and analysis

Basic Request 1: Format: standard_format Process Hide Processings

Backend Information:
05861e73-992a-4363-911b-5ac1d87beeb4 > Logs

REANA backend information can retry jobs on failure

Results Upload to RECAST

- requestwflow-checkmate
- stops_herwig_nllfast-atlas_analysis
- stops_herwig_nllfast-checkmate (0.03|0.09)
- madgraph_pythia-checkmate
- requestwflow-atlas_analysis
- madgraph_pythia-atlas_analysis

CLs obs|CLs exp result for this implementation

RECAST

Live Monitoring Workflow execution

Once, submitted, web-based workflow monitoring — live streaming visualization and logging. Live node-level logging via LogStash (indexable via ElasticSeach for anomaly detection) — similar feel to Travis / GitLab CI job monitoring.

The screenshot shows the RECAST Control Center interface for a job monitor. The browser address bar shows "https://recast-cont...". The page title is "Job Monitor for Workflow 6f35263c-d581-407b-9d4f-1604ec0d3641". The status is "herwig State: RUNNING". The workflow visualization shows a sequence of nodes: "init downstream" leads to "checkmate", which leads to "format_results". The "herwig" node is highlighted in yellow. The log shows messages from the request processor, including "running analysis on worker: worker-3049935782-5cfsv hello".

The screenshot shows the RECAST Control Center interface for a job monitor, displaying detailed log output. The browser address bar shows "https://recast-cont...". The page title is "Job Monitor for Workflow 6f35263c-d581-407b-9d4f-1604ec0d3641". The status is "herwig State: RUNNING". The log output shows progress bars for various stages:

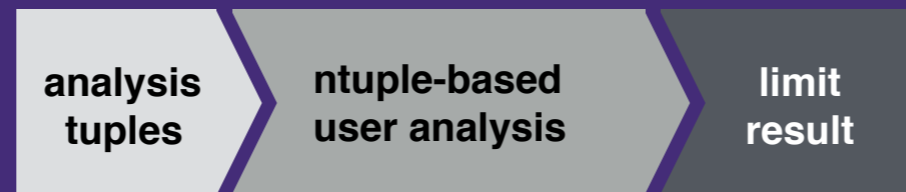
- Progress bar 1: [#####] (93.77%)
- Progress bar 2: [#####] (95.32%)
- Progress bar 3: [#####] (96.90%)
- Progress bar 4: [#####] (98.51%)
- Progress bar 5: [#####] (100.00%)

 The log also includes timestamps (e.g., 2017-08-14T01:59:...) and configuration details such as "Eventfiles: 1", "Analyses: atlas_conf_2013_024 (ATLAS, 0 leptons + 6 (2 b-)jets + etmiss)", and "Output Directories: /srv/workerdata/workdirs/05861e73-992a-4363-911b-5ac1d87beeb4/downstream/checkmate/checkmaterun/r...". The log concludes with "Test: Calculation of CLs from determined signal" and "Warning: Error is dominated by Monte Carlo statistics!".



Within ATLAS we're working towards increasing implementation

- multiple examples of containerized workflows within SUSY and Exotics groups captured via Continuous Integration
- Using REANA prototype deployment to run validation and reinterpretation of these analyses. Cluster Size: $O(1k)$ VCPUs
 - sufficient scale to run large number ntuple \rightarrow limit reinterpretations per day.



- in principle full chain (incl simreco) scan possible on $O(\text{week})$
- proof-of-concept combined upstream MC generation / downstream analysis workflow done



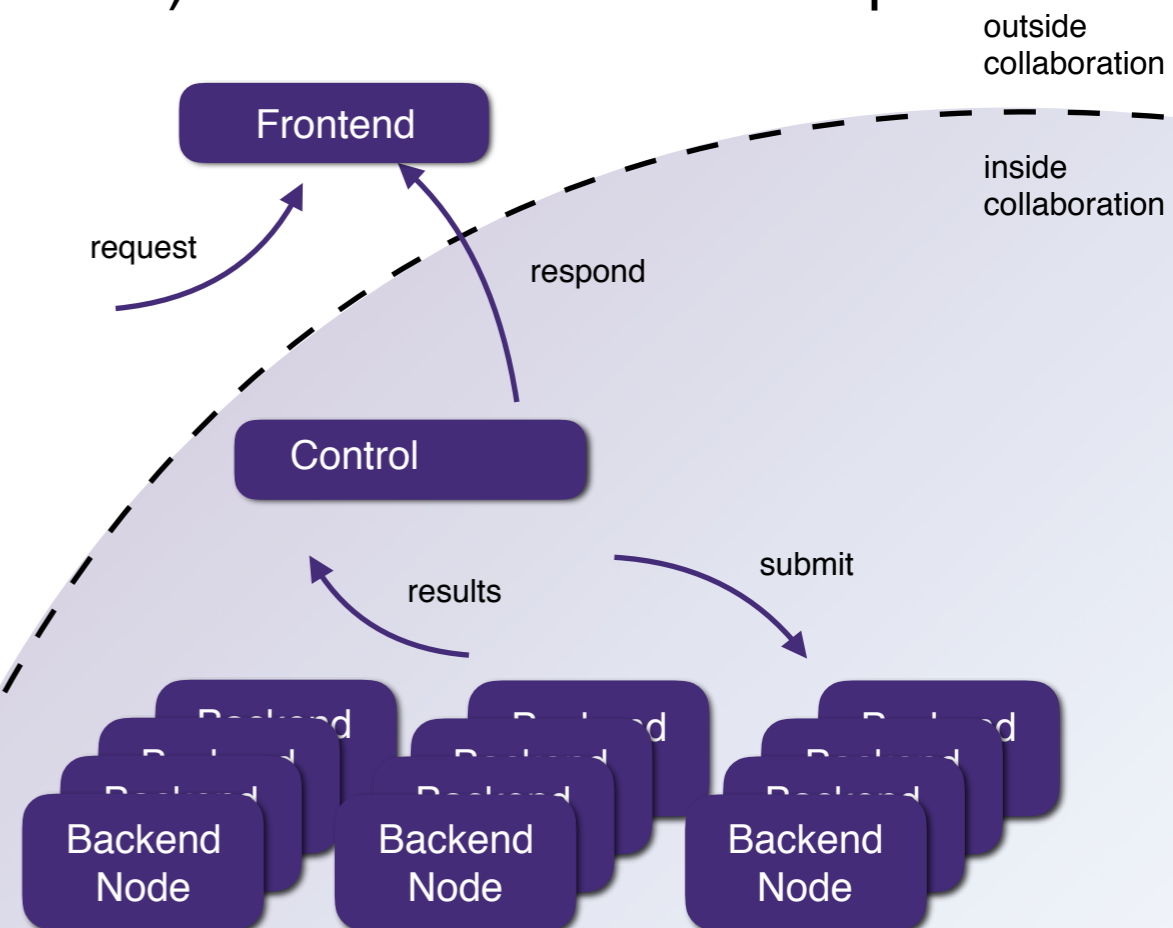
- for now request and processing privileges only ATLAS internal



Future Opportunities:

When implementation of analysis preservation in ATLAS advances and if it proves robust, there is an opportunity to open request privileges to pheno-community

- reinterpetations-as-a-service
- allow scan definition / suggestion by people outside of collaboration, provide e.g. parameter cards based on pre-selection based on external constraints (DM relic density, Higgs Mass, etc..)
- fast-track approval (analysis was already approved) and amendment of HepData records
- Track citation of request for theorist, of response for collaborations (citable data)



Future

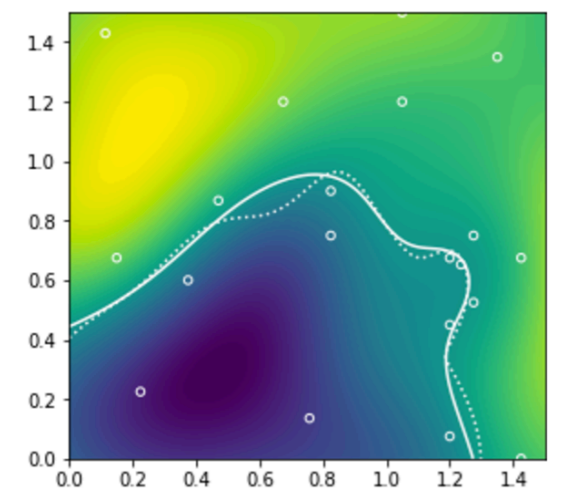
Smart Parameter Scans via Bayesian Optimization

reinterpretation of analysis wrt. to parametrized BSM model = the definition of an expensive multi-variate function. Cartesian grids suffer from curse of dimensionality, need to choose points wisely to find expulsion contour with minimum number of points

$$CL_s = f(m_1, m_2, g_2, g_2, \dots | \mathcal{D}_{13\text{TeV}})$$

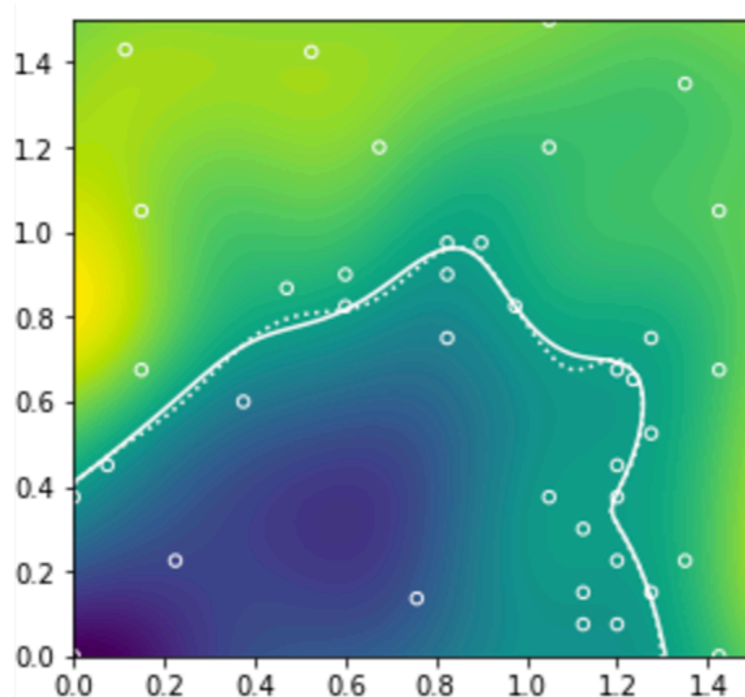
Idea: Bayesian Optimization with objective to minimize overall differential entropy of point classification p.d.f. (i.e. certainty about contour)

Can find much better contour interpolation with **fewer** (*but iterative*) simulated points.

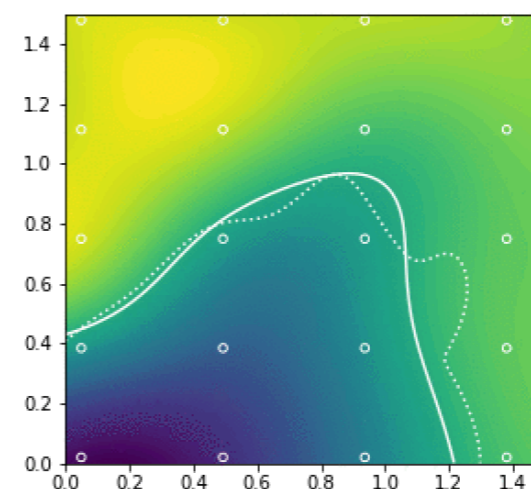


20 points (4 random+ 16 BO)
basic shape already well-known

[animated gif](#)



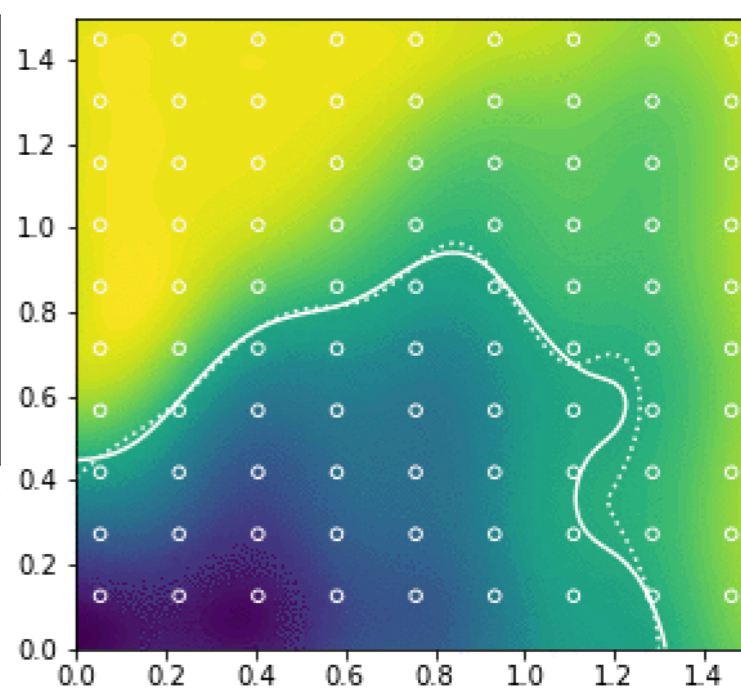
40 points
(mostly detailed refinement)



20 points (5x4) grid

even with dense grid
key contour features
may be missed

[animated gif](#)



90 points in 9x10 grid



- **ATLAS analysis preservation focuses on analysis re-usability via containerized declaratively defined analysis workflows**
- **push to minimize workload on analysis team to capture analysis**
 - **capture during development**
 - **capture only code analysis-related workflow**
 - **easy human-readable text-based workflow capture**
- **Close collaboration with CAP and REANA efforts to drive analysis re-use applications**
 - **native cloud-based solution on modern infrastructure.**
- **First application: collaboration internal re-interpretation of SUSY and Exotics analyses**
- **Exciting Opportunities on the horizon for more systematic approach to reinterpretation**

