



**UNIVERSITÀ  
DI TORINO**

# A Fault Tolerance mechanism for Hybrid Scientific Workflows

**Alberto Mulone**

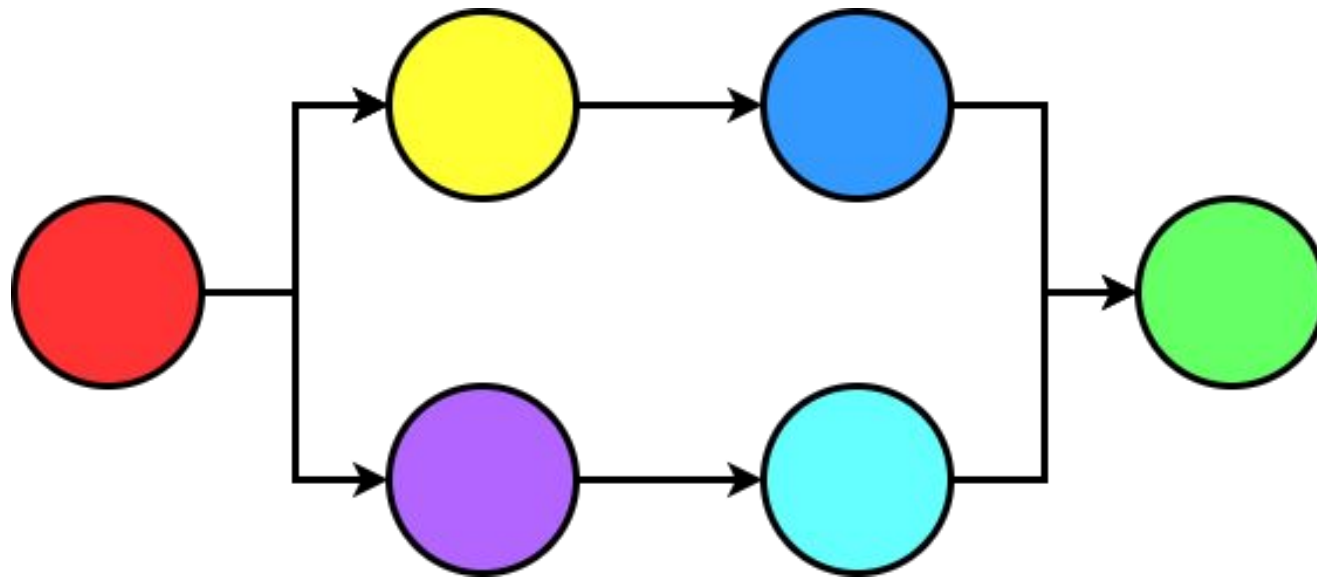
**Doriana Medić**

**Marco Aldinucci**

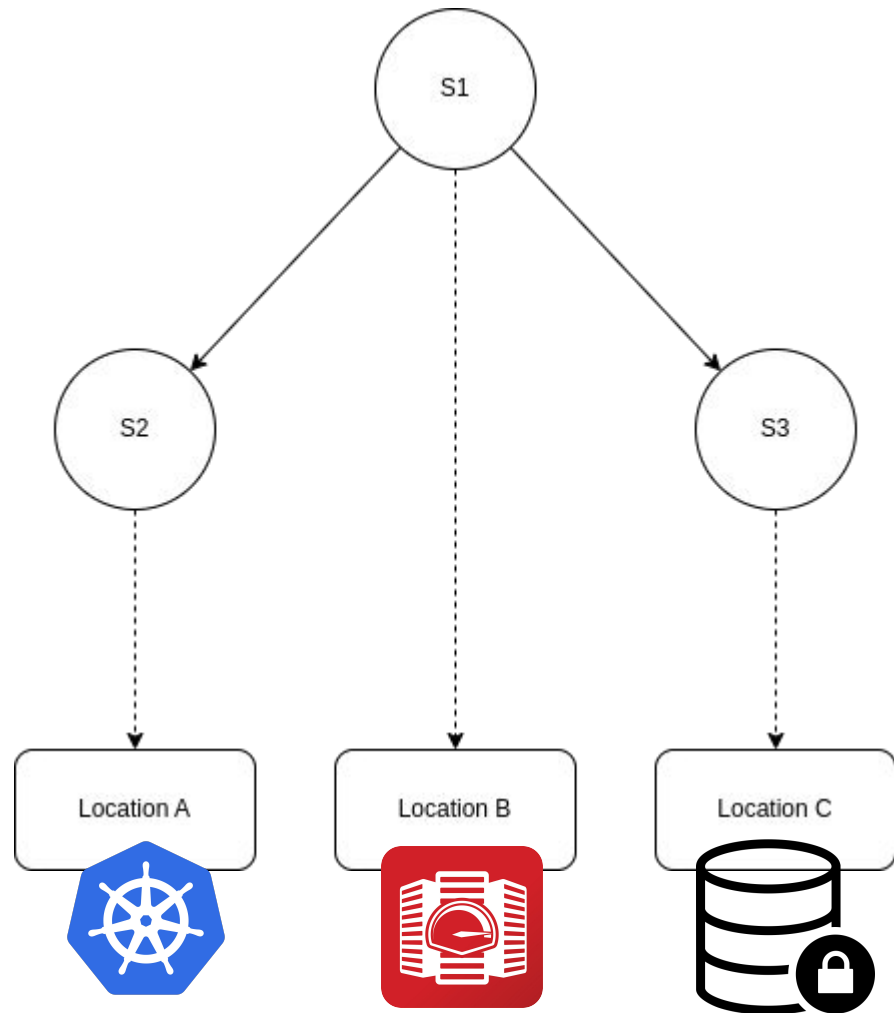
Alberto Mulone, PhD student, University of Turin  
HiPES - EuroPar, Madrid, 26 August 2024

# Workflow

A workflow is an abstraction that models a complex and modular working process as a set of steps and their inter-dependencies.

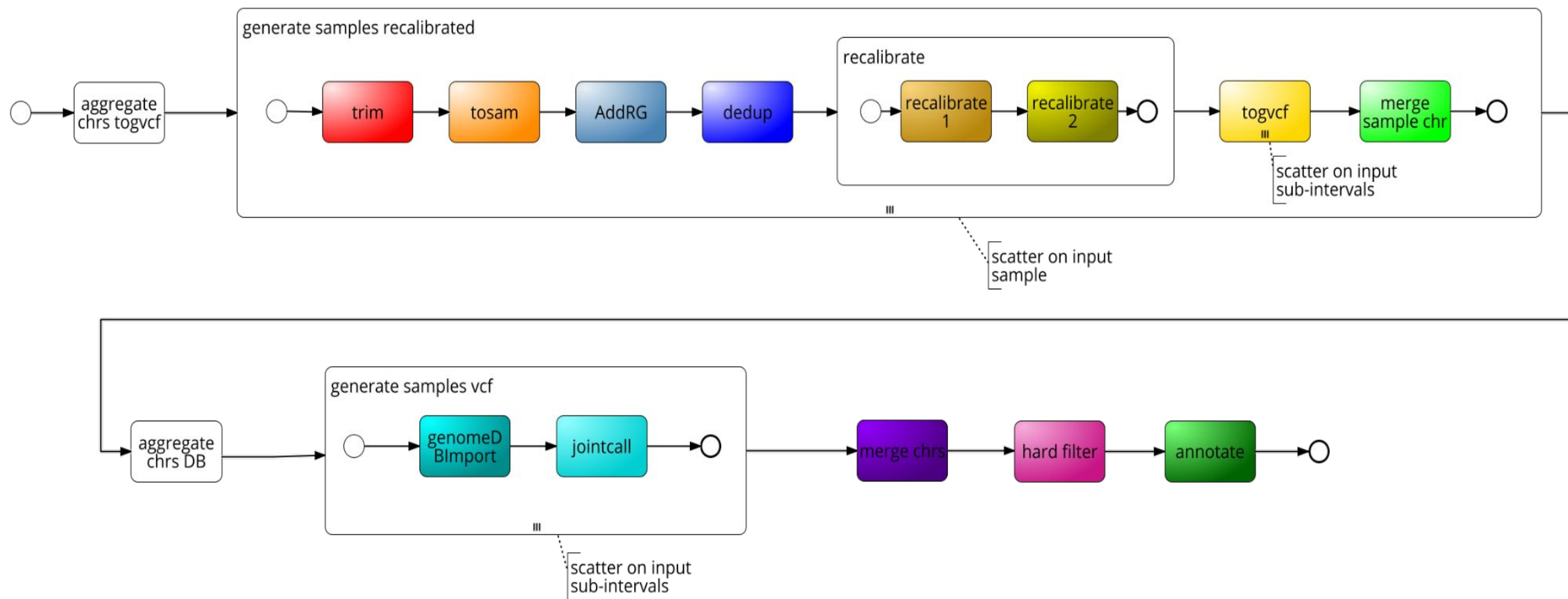


# Hybrid workflows



A hybrid workflow is a workflow whose steps can span multiple, **heterogeneous**, and **independent** computing infrastructures.

# Use case: Variant Calling Pipeline



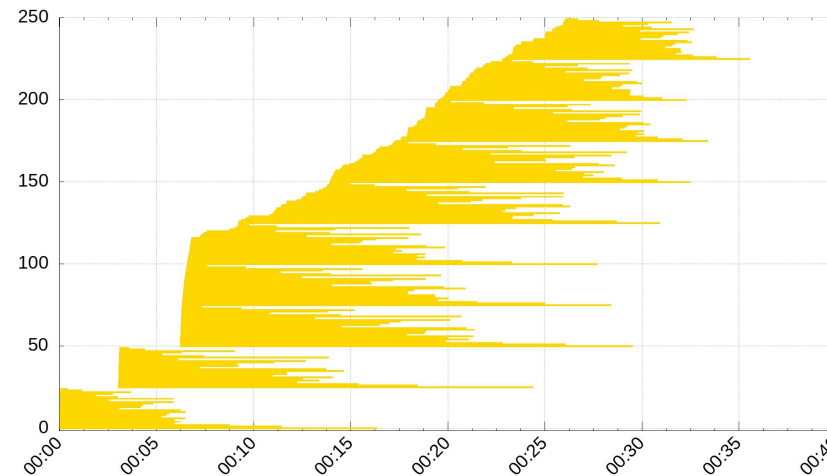
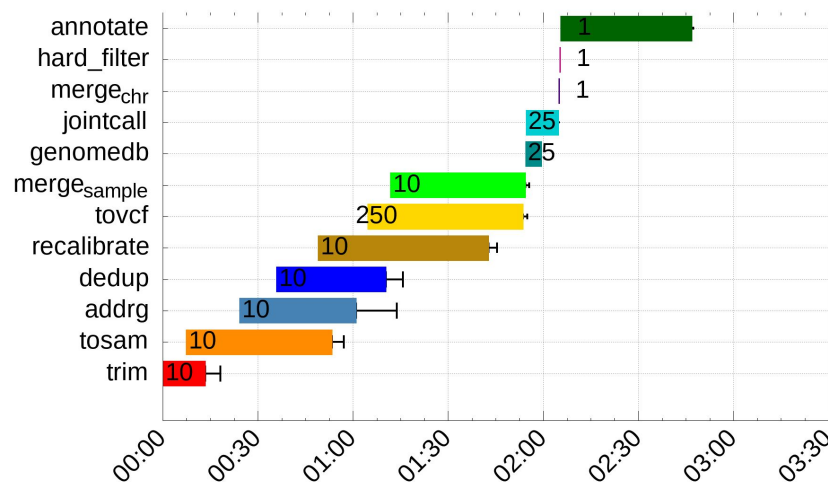
A. Mulone, S. Awad, D. Chiarugi and M. Aldinucci,  
"Porting the Variant Calling Pipeline for NGS data in cloud-HPC environment,"  
2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC),  
Torino, Italy, 2023, pp. 1858-1863, doi: [10.1109/COMPSAC57700.2023.00288](https://doi.org/10.1109/COMPSAC57700.2023.00288).



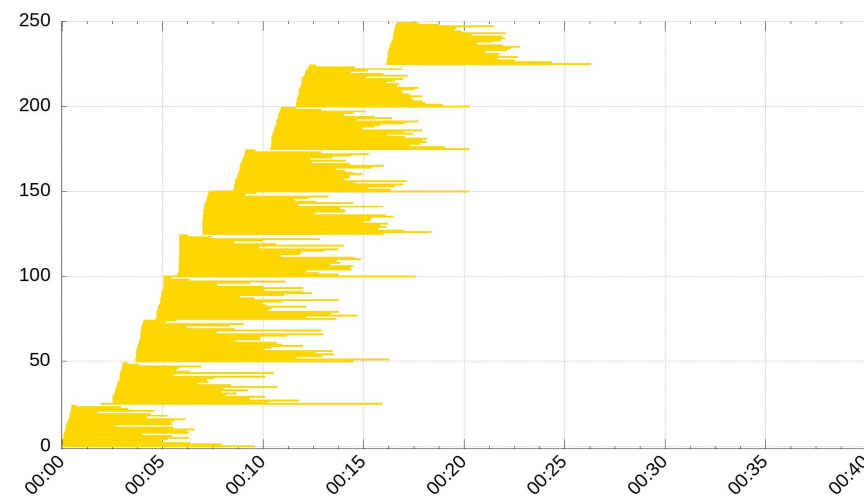
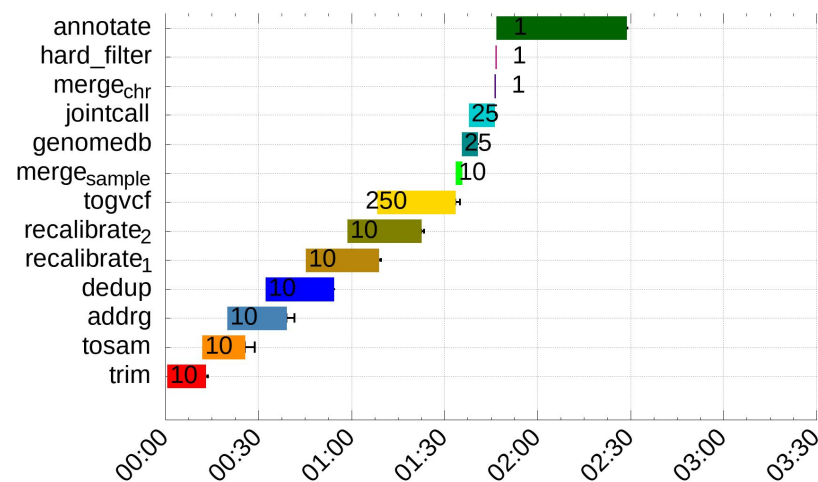


# Use case: Variant Calling Pipeline

VM cloud  
(96 cores)

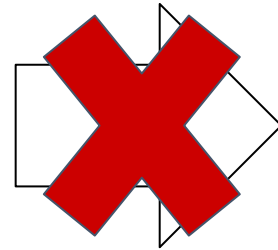
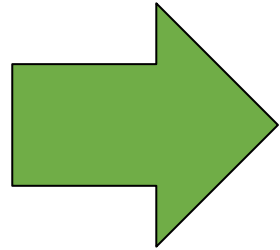


Hybrid  
execution  
(cloud+HPC)



# Failures: soft error

Input data



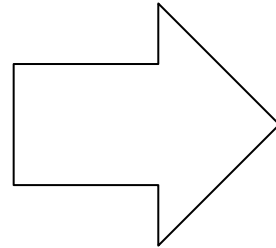
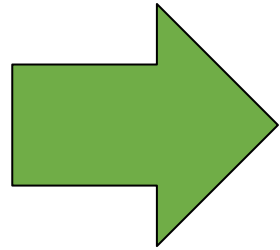
Output data



Application raises  
an exception

# Failures: fail-stop error

Input data



Output data



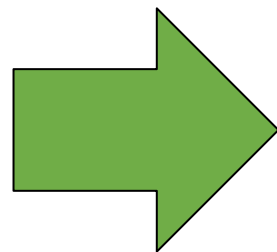
Facility shutdown

# Failures: fail-stop error

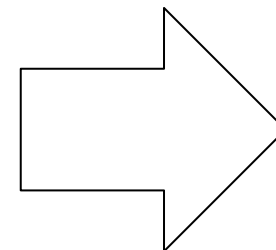
Input data



Data lost



Facility shutdown



Output data



Ephemeral volume

Two blue arrows originate from the text 'Ephemeral volume'. One arrow points upwards and to the left towards the 'Data lost' text. The other arrow points upwards and to the right towards the 'Facility shutdown' text.





# Fault tolerance

A workflow manager is fault tolerant when it can continue the execution correctly even if a failure occurs.

## Our contribution

{ if input data are available, then **retry** failed step  
otherwise **rollback** of steps which output data are lost

Recovery of the failure is delegated to a new sub-workflow. This sub-workflow is called recovery-workflow.

A synchronization mechanism is necessary to manage dependencies across concurrent recovery-workflows.

# Formalization: grammar

## Workflow syntax

$$W ::= \langle l, D_l, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1.t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$
$$\mu ::= \text{exec}(s, I, 0) \parallel \text{tran}(v, l_2) \parallel \underline{\text{tran}(v, l_1)} \parallel \text{rec}(x)$$
$$v ::= d \parallel m_s \parallel m_{d,l} \parallel \text{ok}_{m_s} \parallel \text{ok}_d \parallel \text{err}(x)$$
$$x ::= s \parallel D, l \parallel d, l \parallel m_s, l$$



# Formalization: grammar

Workflow syntax

$$W ::= \langle \mathbf{1}, D_1, t_d \rangle \parallel (W_1, W_2)$$



# Formalization: grammar

Workflow syntax

$$W ::= \langle \mathbf{1}, D_1, t_d \rangle \parallel (W_1, W_2)$$

where  $D_1$  is

$$D_1 = \{ t_w, M(s_i), d, m_{s_i} \}$$

- $t_w$  represents all traces in the workflow
- $M(s_i)$  is the mapping steps-locations
- $d$  is the dataset
- $m_{s_i}$  are the messages that driver sends to locations



# Formalization: grammar

Workflow syntax

$$W ::= \langle 1, D_1, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1 \cdot t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$





# Formalization: grammar

## Workflow syntax

$$W ::= \langle I, D_1, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1 \cdot t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$

- Actions:  $\mu$
- Operators: “.”, “|”, “+”
- Empty trace  $\emptyset$
- Pointer denoting which is the current execution of the trace  $t$ :  $\triangleright t$



# Formalization: grammar

Workflow syntax

$$W ::= \langle \mathbf{1}, D_1, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1 \cdot t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$
$$\mu ::= \mathbf{exec}(s, \mathbf{I}, \mathbf{0}) \parallel \mathbf{tran}(v, \mathbf{1}_2) \parallel \underline{\mathbf{tran}}(v, \mathbf{1}_1) \parallel \mathbf{rec}(x)$$



# Formalization: grammar

## Workflow syntax

$$W ::= \langle I, D_1, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1 \cdot t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$
$$\mu ::= \text{exec}(s, I, 0) \parallel \text{tran}(v, I_2) \parallel \underline{\text{tran}}(v, I_1) \parallel \text{rec}(x)$$
$$v ::= d \parallel m_s \parallel m_{d,1} \parallel \text{ok}_{m_s} \parallel \text{ok}_d \parallel \text{err}(x)$$



# Formalization: grammar

## Workflow syntax

$$W ::= \langle \mathbf{1}, D_1, t_d \rangle \parallel (W_1, W_2)$$
$$t ::= \mu \parallel t_1.t_2 \parallel (t_1 \mid t_2) \parallel (t_1 + t_2) \parallel \triangleright t \parallel \emptyset$$
$$\mu ::= \text{exec}(s, \mathbf{I}, 0) \parallel \text{tran}(v, \mathbf{1}_2) \parallel \underline{\text{tran}}(v, \mathbf{1}_1) \parallel \text{rec}(x)$$
$$v ::= d \parallel m_s \parallel m_{d,1} \parallel \text{ok}_{m_s} \parallel \text{ok}_d \parallel \text{err}(x)$$
$$x ::= s \parallel D, \mathbf{1} \parallel d, \mathbf{1} \parallel m_s, \mathbf{1}$$

# Formalization: rule actions

$$\begin{array}{l}
 \text{(INIT)} \quad \frac{l = \mathcal{M}(s)}{\langle l_d, D_{l_d}, \triangleright \text{init}(t(s), l).t_d \rangle \mid \langle l, D_l, t \rangle \rightarrow \langle l_d, D_{l_d}, \text{init}(t(s), l). \triangleright t_d \rangle \mid \langle l, D_l, t \mid \triangleright t(s) \rangle} \\
 \\
 \text{(EXEC)} \quad \frac{\begin{cases} \text{if } I(s) \not\subset D_l \wedge m_s \in D_l, X = \{\text{err}(D, l)\} \text{ where } D = \{d \mid d \in I(s) \wedge d \notin D_l\} \\ \text{if } I(s) \subset D_l X = \{O(s), M_d\} \text{ where } M_d = \{m_{d,l} \mid d \in O(s)\} \vee X = \{\text{err}(s)\} \end{cases}}{\langle l, D_l, \triangleright \text{exec}(s, I, O).t \rangle \rightarrow \langle l, (D_l \setminus \{m_s\}) \cup X, \text{exec}(s, I, O). \triangleright t \rangle} \\
 \\
 \text{(TRANS)} \quad \frac{v \in D_{l_1} \quad \wedge \quad \begin{cases} \text{if } v \in \{d, m_s\}, \text{ then } X = \emptyset \text{ and } V = \{v, \text{ok}_v\} \vee V = \{\text{err}(v, l_2)\} \\ \text{if } v = \{m_{d,l_1}, \text{ok}_d, \text{ok}_{m_s}\}, \text{ then } X = \{v\} \text{ and } V = \{v\} \end{cases}}{\langle l_1, D_{l_1}, \triangleright \text{tran}(v, l_2).t_1 \rangle \mid \langle l_2, D_{l_2}, \triangleright \text{tran}(v, l_1).t_2 \rangle \rightarrow \\ \langle l_1, D_{l_1} \setminus X, \text{tran}(v, l_2). \triangleright t_1 \rangle \mid \langle l_2, D_{l_2} \cup V, \text{tran}(v, l_1). \triangleright t_2 \rangle} \\
 \\
 \text{(TRANSEERR)} \quad \frac{v = \text{err}(x) \quad \wedge \quad \text{err}(x) \in D_l \quad \text{where } x = \{s, (D, l), (d, l), (m_s, l)\}}{\langle l, D_l, \triangleright \text{tran}(v, l_d).t \rangle \mid \langle l_d, D_{l_d}, \triangleright \text{tran}(v, l).t_d \rangle \rightarrow \\ \langle l, D_l \setminus \{v\}, \text{tran}(v, l_d). \triangleright t \rangle \mid \langle l_d, D_{l_d} \cup \{v\}, \text{tran}(v, l). \triangleright t_d \mid \triangleright \text{rec}(x) \rangle} \\
 \\
 \text{(LOSTDATA)} \quad \frac{\text{err}(d, l) \in D_{l_d} \quad \wedge \quad m_{d_i, l} \in D_{l_d}}{D_{l_d} = D_{l_d} \setminus \{m_{d_i, l}\}}
 \end{array}$$



# Formalization: recovery rule actions

$$(REC(S)) \quad \frac{err(s) \in D_{l_d}}{\langle l_d, D_{l_d}, t_d \mid \mathcal{P} \mathbf{rec}(s) \rangle \mid \langle l, D_l, t(s) \rangle \rightarrow \langle l_d, D_{l_d} \setminus \{err(s)\}, t_d \mid \mathbf{rec}(s) \mathcal{P} \mid \mathcal{P} t_d(s \setminus t_{d_{I(s)}}) \rangle \mid \langle l, D_l, \mathbf{0} \mid t_s \rangle}$$

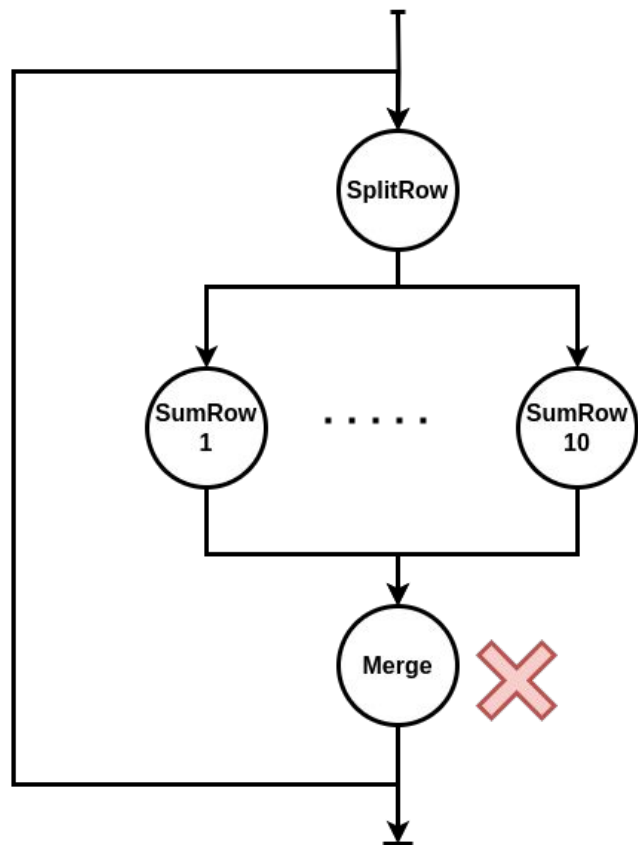
$$(REC(D)) \quad \frac{\forall d_i \in D \text{ such that } err(d_i, l), m_{d_i, l_j} \in D_{l_d} \quad \forall d_h \in D \text{ such that } err(d_h, l) \in D_{l_d} \text{ and } m_{d_h, l'}, d_h \notin D_{l_d} \quad d_h \in I(s_r)}{\langle l_d, D_{l_d}, t_d \mid \mathcal{P} \mathbf{rec}(D, l) \rangle \mid \prod_j \langle l_j, D_{l_j}, t_j \rangle \mid \prod_k \langle l_k, D_{l_k}, t_k \rangle \mid \langle l, D_l, t \rangle \rightarrow \langle l_d, D_{l_d} \setminus \{err(D, l)\}, t_d \mid \mathbf{rec}(D, l) \mathcal{P} \mid \mathcal{P} t_d(s) \mid \mathbf{rec}(s_r) \rangle \mid \prod_j \langle l_j, D_{l_j}, t_j \mid \mathcal{P} \mathbf{tran}(d_i, l) \rangle \mid \langle l, \emptyset, \mathbf{0} \mid t(s) \rangle}$$

$$(REC(D)) \quad \frac{err(d, l), m_{d, l'} \in D_{l_d} \quad \wedge \quad l' \in \mathcal{M}(s_1) \quad d \in O(s_1)}{\langle l_d, D_{l_d}, t_d \mid \mathcal{P} \mathbf{rec}(d, l) \rangle \mid \langle l', D_{l'}, t' \rangle \mid \langle l, D_l, t \rangle \rightarrow \langle l_d, D_{l_d} \setminus \{err(d, l)\}, t_d \mid \mathbf{rec}(d, l) \mathcal{P} \rangle \mid \langle l', D_{l'}, t' \mid \mathcal{P} \mathbf{tran}(d, l) \rangle \mid \langle l, D_l, t \mid \mathbf{tran}(d, l').\mathbf{conf}(ok_d, l) \rangle}$$

$$(REC(m_s)) \quad \frac{err(m_s, l) \in D_{l_d}}{\langle l_d, D_{l_d}, t_d \mid \mathcal{P} \mathbf{rec}(m_s, l) \rangle \mid \langle l, D_l, t \rangle \rightarrow \langle l_d, D_{l_d} \setminus \{err(m_s, l)\}, t_d \mid \mathbf{rec}(d, l) \mathcal{P} \mathbf{tran}(m_s, l).\mathbf{conf}(ok_{m_s}, l) \rangle \mid \langle l, D_l, t \mid \mathbf{tran}(m_s, l_d).\mathbf{conf}(ok_{m_s}, l) \rangle}$$

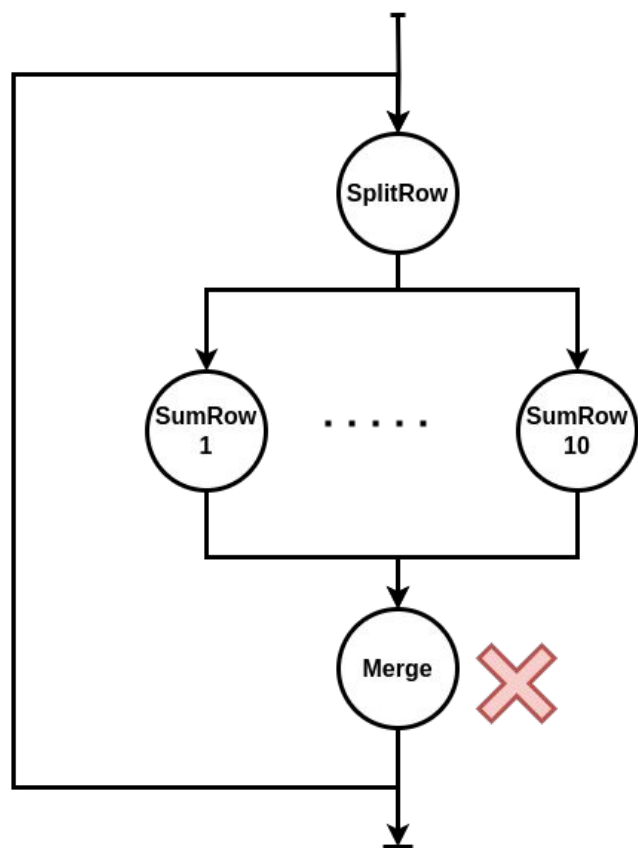
# Example

Workflow

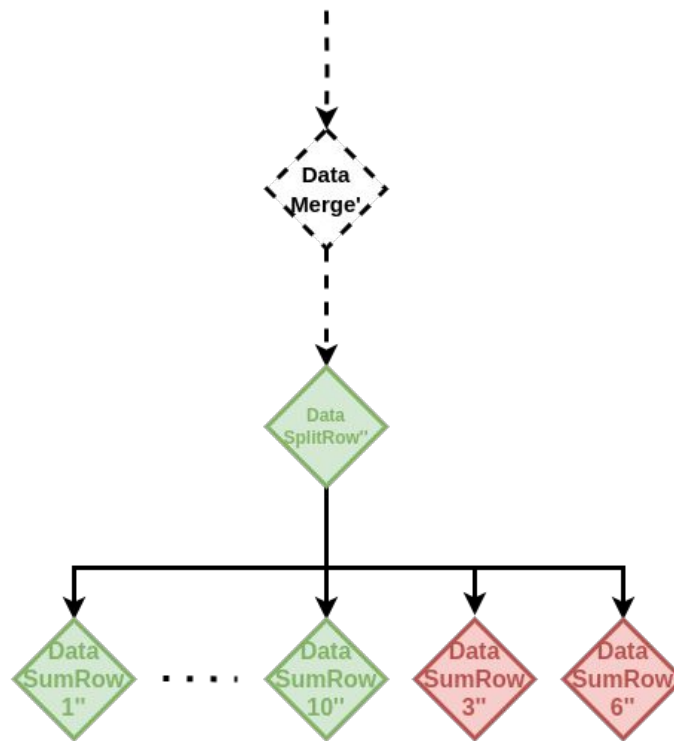


# Example

Workflow

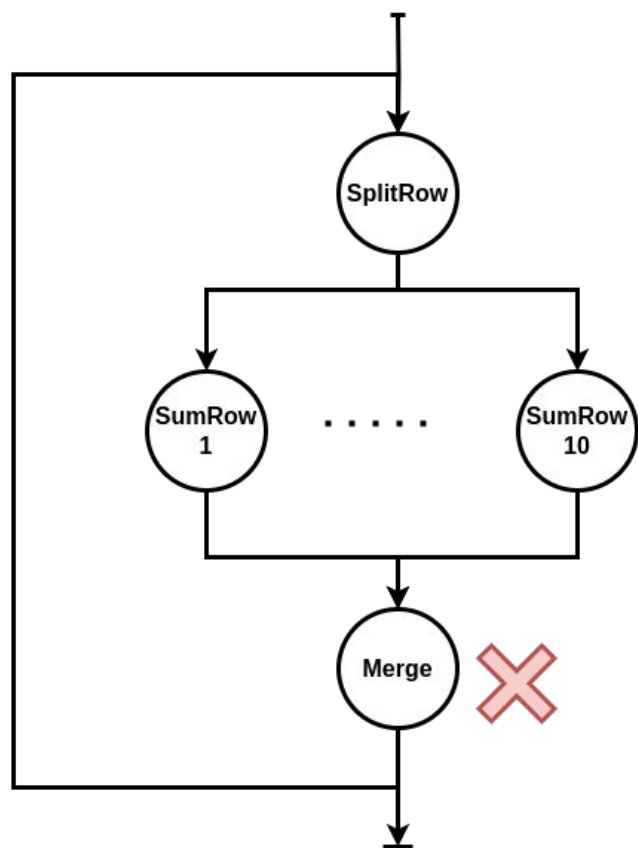


Provenance graph

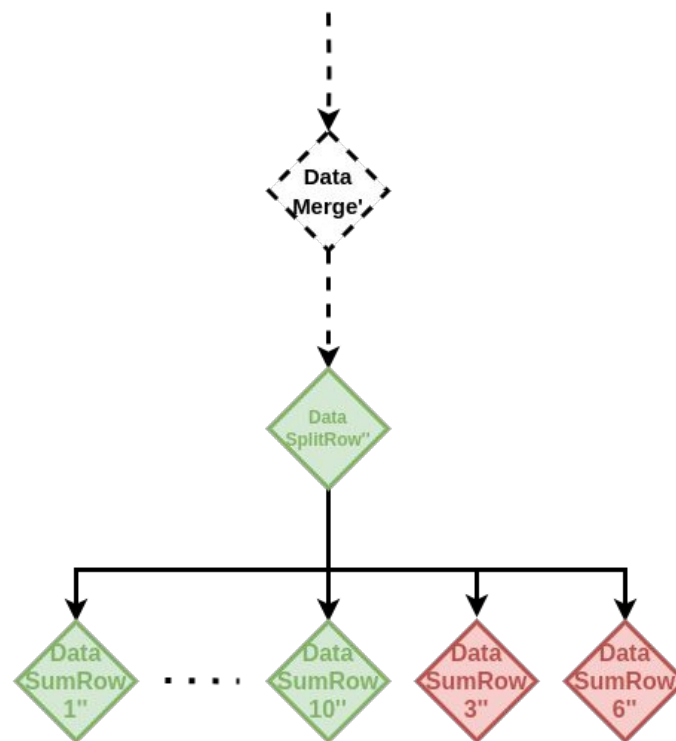


# Example

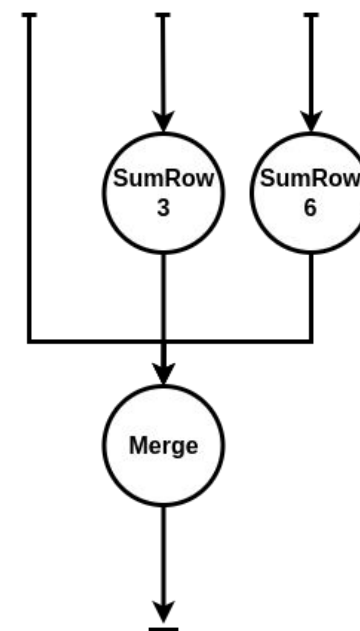
Workflow



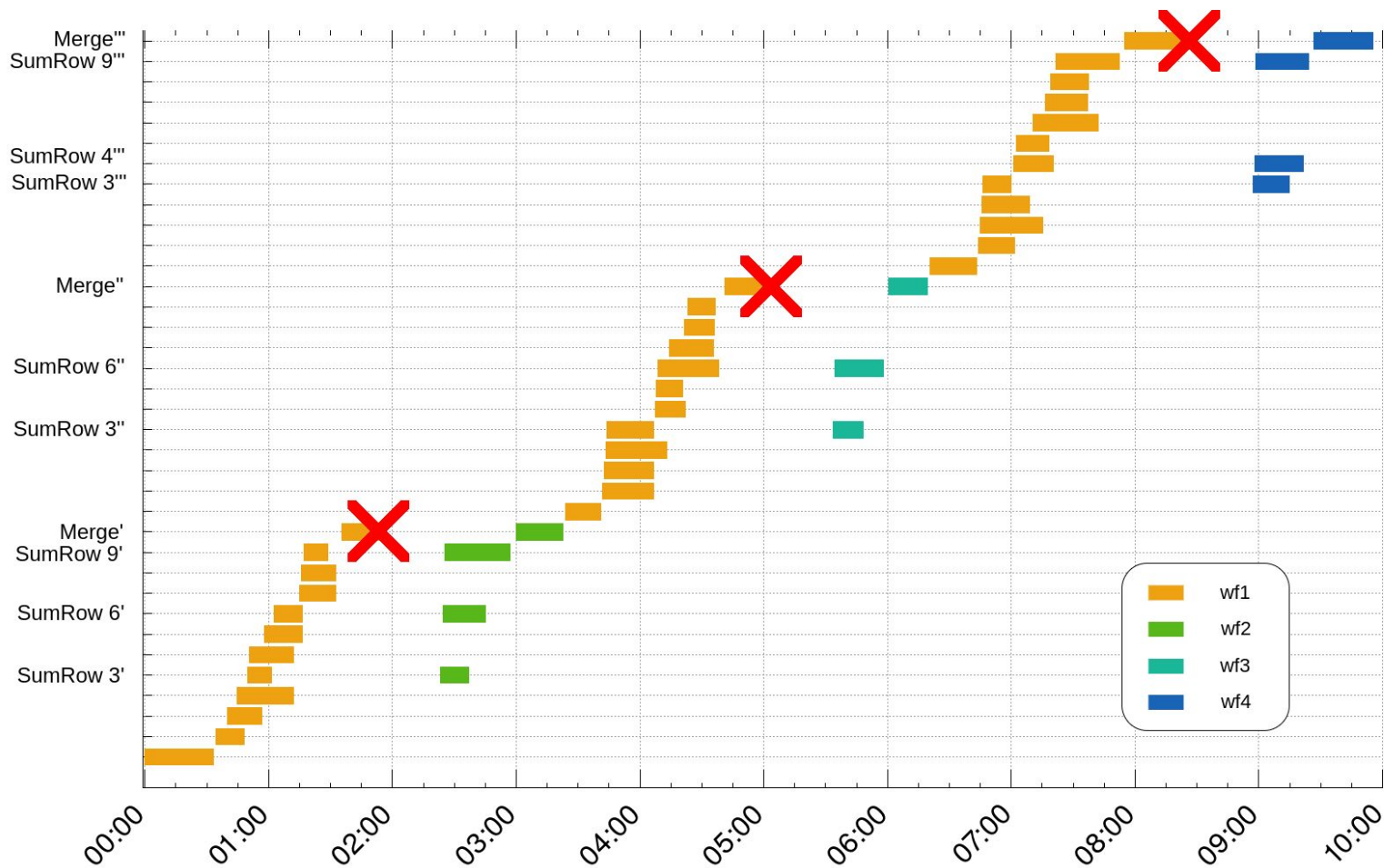
Provenance graph



Recovery-workflow




# Experiment



- Fault tolerance mechanism implemented in StreamFlow (a hybrid WMS)
- Kubernetes with 4 workers
- Injected fail-stop errors on Merge step



# Conclusions

- ↗ Mixed different fault tolerance mechanisms
- $\pi$  Introduced a semantics for fault tolerant hybrid workflows
-  Implemented the mechanism in StreamFlow, a hybrid workflow system

# Future works



Deep performance study about the implemented mechanism in StreamFlow



Extend the semantics with workflow loop



Support non-deterministic workflows