



# Selective and Recurring Re-computation of Big Data Analytics Tasks: Insights from a Genomics Case Study <sup>☆</sup>

Jacek Cała <sup>\*</sup>, Paolo Missier

School of Computing, Newcastle University, Newcastle upon Tyne, UK

## ARTICLE INFO

### Article history:

Received 18 November 2017  
Received in revised form 7 June 2018  
Accepted 20 June 2018  
Available online 14 August 2018

### Keywords:

Re-computation  
Knowledge decay  
Big data analysis  
Genomics

## ABSTRACT

The value of knowledge assets generated by analytics processes using Data Science techniques tends to decay over time, as a consequence of changes in the elements the process depends on: external data sources, libraries, and system dependencies. For large-scale problems, refreshing those outcomes through greedy re-computation is both expensive and inefficient, as some changes have limited impact. In this paper we address the problem of refreshing past process outcomes *selectively*, that is, by trying to identify the subset of outcomes that will have been affected by a change, and by only re-executing fragments of the original process. We propose a technical approach to address the selective re-computation problem by combining multiple techniques, and present an extensive experimental study in Genomics, namely variant calling and their clinical interpretation, to show its effectiveness. In this case study, we are able to decrease the number of required re-computations on a cohort of individuals from 495 (blind) down to 71, and that we can reduce runtime by at least 60% relative to the naïve blind approach, and in some cases by 90%. Starting from this experience, we then propose a blueprint for a generic re-computation meta-process that makes use of process history metadata to make informed decisions about selective re-computations in reaction to a variety of changes in the data.

© 2018 Elsevier Inc. All rights reserved.

## 1. Introduction

In Data Science applications, the insights generated by resource-intensive data analytics processes may become outdated as a consequence of changes in any of the elements involved in the process. Changes that cause instability include updates to reference data sources, to software libraries, and changes to system dependencies, as well as to the structure of the process itself. We address the problem of efficiently restoring the currency of analytics outcomes in the presence of instability. This involves a trade-off between the recurring cost of process update and re-execution in the presence of changes on one side, and the diminishing value of its obsolete outcomes, on the other. Addressing the problem therefore requires knowledge of the impact of a change, that is, to which extent the change invalidates the analysis, as well as of the cost involved in upgrading the process and running the analysis again. Additionally, it may be possible to optimise the re-analysis given prior outcomes and detailed knowledge of, and control over, the analysis process.

### 1.1. Motivation: genomics data processing

In this paper we focus specifically on Genomics data processing, as it is a relevant and paradigmatic case study for experimenting with general re-computation strategies. Next Generation Sequencing (NGS) pipelines are increasingly employed to analyse individuals' exomes (the coding region of genes, representing about 1% of the genome), and more recently whole genomes, to extract insight into suspected genetic diseases, or to establish genetic risk factors associated with some of the most severe human diseases [1–3]. NGS pipelines provide an ideal testbed to study the re-computation problem, as they are relatively unstable and are used to process large cohorts of individual cases. They are also resource-intensive: exome files are of the order of 10 GB each, and a batch of 20–40 exomes is required for the results to be significant. Each 1TB+ input batch requires over 100 CPU-hours to process. Specific performance figures for our own pipeline implementation, which runs on the Azure cloud, can be found in [4].

While the cost and execution time associated to a single execution of these pipelines is decreasing over time [5,4], recent advances in preventive and personalised medicine [6] translate into ambitious plans to deploy genomics analysis at population scale. At the same time, although relatively stable *best practices* are available

<sup>☆</sup> This article belongs to Special Issue: Medical Data Analytics.

<sup>\*</sup> Corresponding author.

E-mail addresses: [Jacek.Cala@ncl.ac.uk](mailto:Jacek.Cala@ncl.ac.uk) (J. Cała), [Paolo.Missier@ncl.ac.uk](mailto:Paolo.Missier@ncl.ac.uk) (P. Missier).

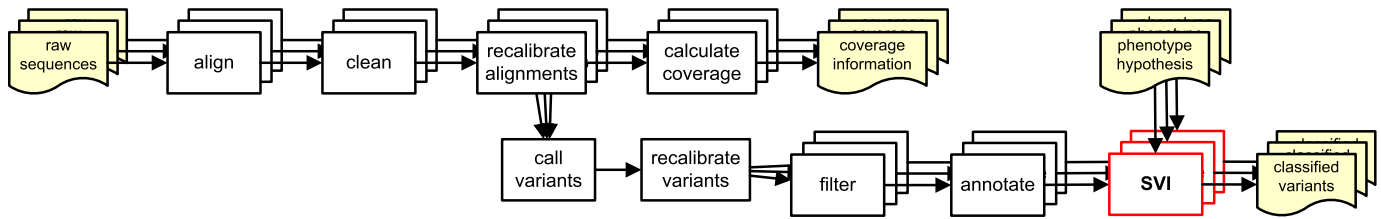


Fig. 1. The Next Generation Sequencing pipeline; highlighted is the variant classification step.

to describe the general structure of the analysis process,<sup>1</sup> their implementations make use of algorithms and tools that are subject to frequent new releases, as well as of reference databases that undergo regular revisions.

In this setting, failing to react to important changes results in missed opportunities to improve on an individual's genetic diagnosis. On the other hand, over-reacting to each and every change is impractical and inefficient, as in many cases the benefits of refresh may be marginal. Using genomics data processing as a case study, we are therefore motivated to explore techniques for *selective* and *incremental re-computation* that optimise the use of the available computing resources vis-à-vis the expected benefit of knowledge refresh on a population of prior outcomes.

### 1.2. Reacting to changes: a meta-process

To clarify the meaning of *selectivity* and *incremental re-computation* in this context, consider: a collection  $C$  of cases, e.g., a cohort of individuals' genomes; an analysis process  $P$ , e.g. an NGS pipeline; a collection of executions of  $P$  on each input  $x_i \in C$ , which generate corresponding outcomes  $y_i$  with processing cost  $c_i$ ; and a set  $D = \{d_1 \dots d_m\}$  of versioned dependencies, i.e., software libraries or reference databases. When a new version  $D'_j$  of a dependency  $D_j \in D$  becomes available, we expect the change  $D_j \rightarrow D'_j$  to have different impact on different outputs  $y_i$  computed at some earlier time: some of these outputs will be unaffected, while others will be partially or completely invalidated, as we will show in examples later.

We are going to define *impact* in terms of a change on a specific output  $y_i$  in terms of some type-specific *diff* functions that compute the differences between two versions  $y_i, y'_i$  of an output. Assuming that expected impact can be estimated, we define the *scope* of the change as the subset of  $C' \subseteq C$  of inputs  $x_i$  such that the change will have non-zero impact on the corresponding output  $y_i$ , and the *selectivity* of the change as  $1 - \frac{|C'|}{|C|}$ . Those  $x_i \in C$  that are within the scope of a change are candidates for re-computation, and it may be possible to prioritise them using knowledge of the cost  $c_i$  of their earlier processing, the quantified extent of impact, along with domain-specific knowledge of their relative importance (for instance, more severe genetic diagnoses). Such considerations, however, are beyond the scope of this paper.

Instead, here we study techniques to (i) estimate the scope of a change, without having to recompute each output, and (ii) perform *incremental* re-computation: given a *white box* specification of  $P$ , for instance as a script or as a workflow, we want to efficiently identify the minimal fragment of  $P$  that is affected by the change, in order to optimise the re-computation of the  $x_i$  that are within the scope of the change. We define such techniques within the framework of the ReComp *meta-process*. ReComp takes as input a history of prior analysis and a change event, as indicated above, and controls the incremental re-execution of the underlying

process  $P$  on selected inputs that are within the scope of the change.

Not all scenarios involving  $C, P$ , and changes in  $P$ 's dependencies are equally suitable for optimisation using ReComp, however. Specifically, ReComp is most effective when changes have high selectivity (only few of the cases are affected), when process  $P$  is a *white box*; and when the change affects only a few of  $P$ 's components, providing scope for incremental re-computation. In the next section we select our target case study following these three requirements, by analysing three scenarios involving different reference data and software tool changes within the realm of Genomics. Firstly, however, we must briefly describe NGS pipelines.

### 1.3. Variant calling and interpretation

Fig. 1 depicts the anatomy of the NGS pipeline implementation available from our lab. It consists of two main phases: (i) exome analysis and variant calling and annotation [4], and (ii) variant interpretation [7]. The first phase closely follows the guidelines issued by the Broad Institute.<sup>2</sup> It takes a batch of raw input exomes and, for each of them, produces a corresponding list of variants, or *mutations*, defined relative to the current reference human genome (in the order of tens of thousands). Particularly critical in this phase are the choices of reference genome, currently at version h19, and the choice and version of the *variant caller*. Currently we use FreeBayes [8], one of several such algorithms [9]. At the end of this phase, each variant will have been annotated using a variety of statistical predictors of the likelihood that the variant contributes to a specific genetic disease.

Only a very small fraction of these variants are deleterious, however. The second phase, which we have called Simple Variant Interpretation (SVI in the figure), aims to identify those the few tens of variants that may be responsible for an individual's phenotype, i.e., the manifestation of a suspected genetic disease. In addition to using the predictors, SVI also makes use of databases that associate phenotype descriptions with sets of genes that are known to be broadly implicated in the phenotypes, such as OMIM GeneMap.<sup>3</sup> It also uses databases of known variants and their deleteriousness such as NCBI ClinVar,<sup>4</sup> HGMD,<sup>5</sup> and possibly others.<sup>6</sup>

In more detail, the SVI portion of the pipeline consists of three main steps (Fig. 2): (1) mapping the user-provided clinical terms that describe a patient's phenotype to a set of relevant genes (*genes-in-scope*), (2) selection of those variants that are in scope, that is, the subset of the patient's variants that are located on the genes-in-scope, and (3) annotation and classification of the variants-in-scope according to their expected pathogenicity. Classification consists of a simple traffic-light system {red, green, and amber} to denote pathogenic, benign and variants of unknown or uncertain pathogenicity, respectively. In this process, the class of a

<sup>2</sup> <https://software.broadinstitute.org/gatk/best-practices>.

<sup>3</sup> <http://data.omim.org>.

<sup>4</sup> <https://www.ncbi.nlm.nih.gov/clinvar>.

<sup>5</sup> <http://www.hgmd.cf.ac.uk>.

<sup>6</sup> [http://grenada.lumc.nl/LSDB\\_list/lsdbs](http://grenada.lumc.nl/LSDB_list/lsdbs).

<sup>1</sup> <https://software.broadinstitute.org/gatk/best-practices>.

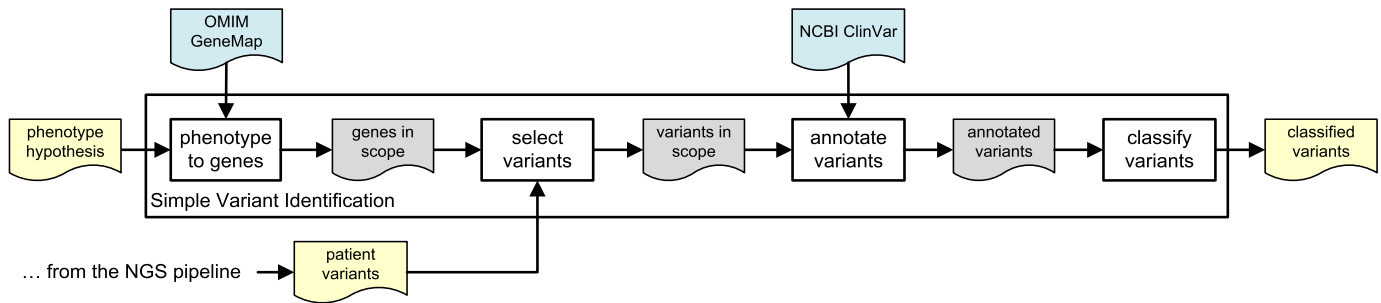


Fig. 2. The high-level architecture of the SVI process.

variant is determined simply by its *pathogenicity status* as reported in ClinVar. Importantly, if any of the patient variants is marked as red, the phenotype hypothesis is deemed to be confirmed, with more red variants interpreted as stronger confirmation.

#### 1.4. Candidate re-computation scenarios

We now present three real scenarios for changes to the processing pipeline just described, as candidates for our in-depth experimentation: (i) a step change in the reference genome assembly, (ii) version upgrade to the variant caller, and (iii) updates to one of the SVI reference databases, ClinVar.

##### 1.4.1. Step change in reference genome assembly

The reference genome is currently undergoing major changes within the bioinformatics community. The history of how the Genome Reference Consortium (GRC) managed the progression of the human genome assembly since 2007 is summarised for instance in [10]. While this provides detailed insight for the interested reader, for our purposes it suffices to note that the Global Alliance for Global Health<sup>7</sup> is working on a new reference genome, h38, that is so drastically different from its predecessors, to require a re-design of most tools and of the entire pipeline. There are two main reasons why h38 will be disruptive. Firstly, it will be graph-structured, taking into account multiple possible transcriptions of the same gene (i.e., during protein synthesis), and secondly, it is the first coordinate-changing assembly update since 2009 [11]. Such disruptive step-changes are rare, however, as the current genome assembly, h19, has been stable for a number of years and is likely to remain in use for quite some time. From the ReComp perspective, this change is likely to have very low selectivity, i.e., every case in  $C$  (every genome ever processed) will be affected, while not leaving much space for fine-grained selection of sub-processed with the established pipeline  $P$ , because most of its elements will be disrupted.

##### 1.4.2. Variant caller version change

Complementary to updates in reference datasets, new releases for one of the tools that make up the pipeline also represent notable change events that may trigger re-computation. The FreeBayes caller we use in our pipeline, for instance, has seen multiple releases between 12/2013 (v0.9.10) and 04/2018 (v1.2.0, current at the time of writing). To assess the broad impact of these changes, we have compared the output variant sets for 16 patients using three versions of the caller, namely v0.9.10, v1.0.2 (12/2015), and v1.1 (11/2016). The results, shown in Fig. 3, are consistent with other, more extensive comparative studies like [12]. In particular, we can see that over 50,000 of the variants that appear in the v0.9.x output are no longer identified as such in v1.0.2, representing a substantial 10.3% false positive detection over the previous

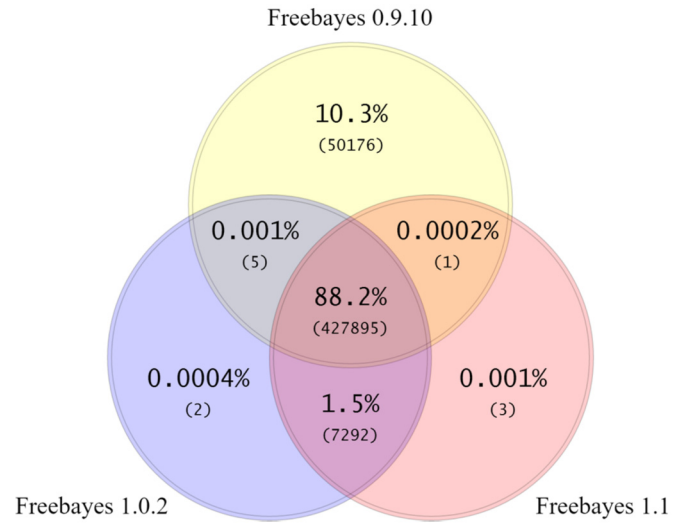


Fig. 3. Variations in variants produced by different versions of the FreeBayes caller.

version. Conversely, the minor version upgrades are much more consistent with each other. This provides empirical evidence of instability of analysis outcomes especially in the early releases of new critical algorithms as part of established pipelines.

##### 1.4.3. Updates to the SVI reference databases

Our third candidate change scenario involves version changes in ClinVar, one of the reference databases used in the SVI portion of the pipeline. We analysed the variants for a cohort of 33 patients for three distinct phenotypes: *Alzheimer's disease*, *Frontotemporal Dementia-Amyotrophic Lateral Sclerosis (FTD-ALS)* and the *CADASIL syndrome*. For each patient we ran SVI using consecutive monthly versions of ClinVar, from 07/2015 to 10/2016, for a total of 16 re-runs per patient, and recorded whether the new version would have modified a diagnosis that had been obtained using the previous version. A change in diagnosis occurs when new variants are added to the selection, others are removed, or existing variants change their classification because their status in ClinVar has changed.

Table 1 summarises the results. We recorded four types of outcomes. Firstly, confirming the current diagnosis (■), which happens when additional variants are added to the red class. Secondly, retracting the diagnosis, which may happen (rarely) when all red variants are retracted, denoted ♦. Thirdly, changes in the amber class which do not alter the diagnosis (□), and finally, no change at all (•).

These results, however limited in scope, suggest good selectivity for this type of change. Indeed, the majority of the changes reported here are ultimately of low interest to clinicians, and so greedy re-computation would be highly inefficient. This comes as little surprise because some human genetic diseases tend to be un-

<sup>7</sup> <http://genomicsandhealth.org/>.

**Table 1**

Changes observed in the output of the SVI tool for a cohort of 33 patients following updates in the NCBI ClinVar reference database between July 2015 and October 2016; ■ and ♦ – significant (positive and negative) changes in the SVI output, □ – insignificant change in the output, ‘.’ – no change at all.

Phenotype hypothesis	Frontotemporal Dementia-Amyotrophic Lateral Sclerosis											CADASIL											Alzheimer's disease														
	Variant file	D_1071	D_1049	D_1041	D_0899	D_0854	D_0830	C_0171	C_0098	C_0056	C_0053	C_0051	B_0307	D_1136	C_1457	C_0072	C_0071	C_0068	C_0065	B_0396	B_0384	B_0370	B_0365	B_0358	B_0338	B_0331	B_0229	B_0214	B_0209	B_0208	B_0203	B_0202	B_0201	B_0198			
08/15		■	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
09/15		□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
10/15		.	.	.	.	.	.	.	.	.	.	.	.	■	□	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
11/15		□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□		
12/15		□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
01/16		.	.	.	.	.	.	.	.	.	.	.	.	.	□	□	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
02/16		.	.	□	.	□	.	□	□	□	□	□	□	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
03/16		□	.	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
04/16		.	.	.	.	.	.	.	.	.	.	.	.	.	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
05/16		.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
06/16		.	.	□	.	□	.	□	□	□	□	□	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
07/16		.	.	□	.	□	.	.	□	□	□	.	.	.	□	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
08/16		.	.	.	.	.	.	.	.	.	.	.	.	♦	♦	■	♦	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
09/16		.	.	.	.	.	.	.	.	.	.	.	.	.	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	
10/16		□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	□	

**Table 2**

Comparison between of changes and expected ReComp effectiveness.

Change				ReComp effectiveness	
Scenario	Type	Rate	Granularity	Selectivity	Optimisation opportunities
h19 → h38	Data	Low (years)	Coarse	Low: high impact over entire cohort	Low: entire process disrupted
Variant caller version change	Algorithm	Every few months	N/A	High: most cases only marginally affected	High but straightforward: always restart from variant calling
SVI	Data	Monthly	Fine (individual records)	High: most cases unaffected, many not significantly affected, very few changes in diagnosis	High: partial, differential re-run possible

derpinned by a very few rare variants [13], whilst those associated with common diseases (as above) are widely studied. Therefore, the knowledge about them is quite stable, especially when considered on a monthly time scale. This also suggests that rare diseases may provide a more compelling case for selective re-computation, as knowledge about them is more likely to evolve over time. Finally, we note that some updates have a higher impact than others, for instance the 08/2016 release of ClinVar.

1.4.4. Choice of target experimental study

The characteristics of the changes just presented are summarised in Table 2. As noted above, a reference genome change results in low sensitivity and little chance for optimisation of process re-run, limiting the effectiveness of ReComp in this case. In contrast, both the variant caller version change and the SVI reference data changes are good candidates, providing potentially good selectivity. Compared to a change in software, however, changes in reference data have the additional advantage that we can apply techniques based on fine-grained differencing of the dataset versions, such as those presented in Sec. 5, making this the case study of choice to illustrate ReComp's capabilities. Noting that the changes only affect the SVI portion of the pipeline, our experiments are focused on this final part of the pipeline.

1.5. Paper contributions

Our main and novel contribution is the specification of a *generic selective re-computation meta-process*, which harnesses an underlying Big Data process and seeks to optimise the extent of its re-executions in reaction to each data change, relative to a blind re-computation baseline. Importantly, we ensure that re-computation is *lossless*, i.e. each outcome on which the change has non-zero impact is indeed updated. The meta-process combines four basic steps that we describe semi-formally. Within this context we propose an algorithm to address one of the steps, scope identification, and also observe that processes distributive over set union and difference can effectively perform differential execution. Our *second contribution* is an extensive experimental study, conducted using the SVI process as testbed, to determine the effectiveness of the meta-process and assess its limitations. Finally, the *third contribution* is the outline of the remaining challenges in addressing re-computation and a discussion about the ideas for a more comprehensive technical approach.

2. A generic meta-process for selective re-computation

As mentioned, the meta-process includes four macro steps: (S1) computing differences between old and new datasets that contribute to the underlying process, (S2) determination of the scope

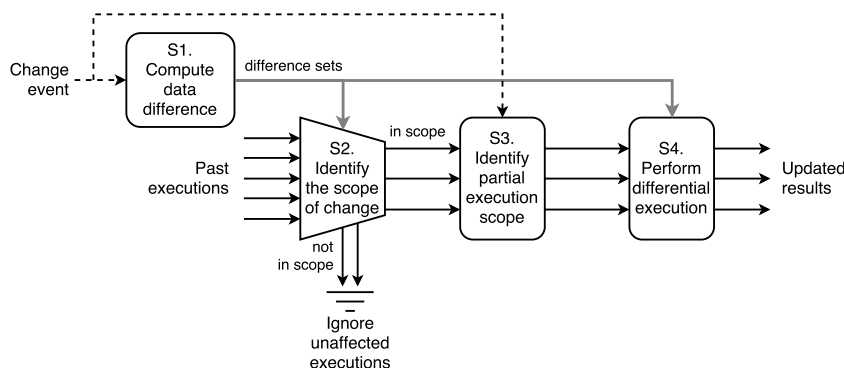


Fig. 4. Overall design of the generic selective re-computation meta-process.

of affected past executions, (S3) identification of the fragment of the underlying process that is affected by the change and (S4) differential execution (Fig. 4).

**S1: Data difference** Data differences are computed in reaction to any data change event observed in the environment, and consist of three sets of records: added, removed and updated. Here the main challenge is to define domain-specific, semantically-rich differencing functions that capture only most relevant changes.

**S2: Scope of change** Scope is defined relative to the population of outcomes from past executions, for instance a large cohort of patients. Here the challenge is to accurately identify all of the past outcomes (patient diagnoses) that are affected by the change. In the example presented in Table 1 a perfect oracle would recognise that only 14 instances of past outcomes need to be considered, namely one instance for each ClinVar release 08/15, 10/15, 11/15 and eleven instances for ClinVar 08/16. In reality, however, it is not always possible to accurately predict which occurrences are affected by the change as the analytics process may implement a complex algorithm (e.g. a simulation of a physical model) or the change itself does not provide enough context to understand its impact. But often the scope can be split into the set of clearly affected and clearly unaffected instances, and so we can use it to reduce the re-computation effort. To address the scope challenge, we rely on process- and data-specific *impact functions* that take data difference and past outcomes, and help eliminate executions not affected by the change.

**S3: Partial re-execution** If we have insight into the structure and semantics of analytics process  $P$ , we may be able to reduce the re-computation of  $P$  to only those parts that are located *downstream* from the point where the changed dataset is first used. For this, we are inspired by techniques for *smart rerun* of workflow-based applications [14,15].

The basic requirement to support partial re-execution is that intermediate data of every past execution be cached, so it can be re-used in lieu of re-executing part of a process that is known to have not been affected by the changes. For example, if the updated data is used only in the middle of the process we may be able to skip processing of its initial part and use previously computed intermediate data instead. This is much easier to achieve in the case of scientific workflows (dataflows) such as our NGS pipeline, where the data dependencies are explicit. In this scenario, the main difficulty is to find a good balance between how much intermediate data should be cached, versus how much we should re-generate in order to minimise the overall re-computation cost (monetary, runtime and/or storage). Some guidance on this topic has been presented e.g. in [16,17]. We present extensive experiments in Sec. 7.

**S4: Differential execution** Finally, the last step to optimising re-computation involves the use of differences between versions of input and reference data. That is especially important in the Big Data analyses, when changes often affect only a small part of large input data. However, whether or not it makes sense to execute a process using a difference rather than complete input data depends on the function it implements.

Differential dataflows [18] is a promising approach to realise differential execution. Depending on the amount of change in the input data, the differential dataflows framework<sup>8</sup> can offer very high performance gains and reduce runtime up to three orders of magnitude. However, to exploit the full potential of differential dataflows in the NGS pipeline all the algorithms used in the pipeline would need to be rewritten following that approach – a non-trivial task, in general. In Sec. 6 we present our experiments and observations involving re-computation of SVI using difference sets.

## 2.1. Notation

Here we provide a simple reference framework for expressing steps S1–S4 and introduce the technical elements that underpin our experiments. Consider an instance  $P_i$  of a deterministic analytics process  $P$ , which takes input  $x_i$  and produces output  $y_i$ , using reference datasets  $D = \{D_1 \dots D_m\}$ .  $D$  is typically the same dataset across a population  $x_1 \dots x_n$  of inputs. For SVI,  $D = \{OM, CV\}$  consists of the two reference databases, OMIM GeneMap and NCBI ClinVar as mentioned earlier.

The  $x_i$  and  $y_i$  may be tuple-valued:  $x_i = \langle x_{i1} \dots x_{in} \rangle$ ,  $y_i = \langle y_{i1} \dots y_{im} \rangle$ . We denote the types of  $x_{ij}$  (resp  $y_{ij}$ ) in each instance  $x_i$  (resp  $y_i$ ) with  $X_{ij}$  (resp  $Y_{ij}$ ). For instance, in SVI, patient  $i$  is represented by input  $x_i = \langle x_{i1}, x_{i2} \rangle$  where  $x_{i1}$  is the list of patient variants and  $x_{i2}$  is a set of *phenotype terms*.<sup>9</sup> The patient's diagnosis is the single output  $y_i = \langle y_{i1} \rangle = \{(v, c)\}$ , a set of variants  $v$  along with their class label  $c$ , i.e., Green, Amber, or Red. For simplicity of notation, and without loss of generality, in the following we are going to refer to inputs and outputs simply as  $x_i$ ,  $y_i$ , and to their types as  $X, Y$ .

**Data versions and change notation.** Each of the  $x_i$  and  $D_j \in D$  may have multiple versions, which change over time. We denote the version of  $x_i$  at time  $t$  as  $x_i^t$ , and the state of  $D_j$  at  $t$  as  $D_j^t$ . We write  $x_i^t \rightarrow x_i^{t'}$  to denote that a new version of  $x_i$  has become available at time  $t'$ , replacing the version  $x_i^t$  that was current at  $t$ . Similarly,  $D_j^t \rightarrow D_j^{t'}$  denotes a new release of  $D_j$  at time  $t'$ .

<sup>8</sup> <https://github.com/frankmcscherry/differential-dataflow>.

<sup>9</sup> Phenotype is a description of patient's disease or condition expressed using terms from a formal vocabulary, such as OMIM or the Human Phenotype Ontology.

**Executions.** We denote the execution  $P_t$  of  $P$  that takes place at time  $t$  by:

$$\langle y_i^t, c_i^t \rangle = \text{exec}(P, x_i^t, D^t) \quad (1)$$

where  $D^t = \{D_1^t \dots D_m^t\}$ .  $c_i^t$  denotes the cost of the execution, for example a time or monetary expression that summarises the cost of cloud resources. We also assume for simplicity that  $P$  remains constant.

**Current outcomes.** Finally, by slight abuse of notation, with  $Y^t = \{y_1^t, y_2^t \dots y_N^t\}$  we denote a set of  $N$  outcomes that are *current* at time  $t$ , i.e., each  $y_i^t$  is the latest in a series of values  $y_i^{t_1} \dots y_i^{t_k}$  with  $t_k \leq t$ .

**Impact of a change.** We say that change  $D_j^t \rightarrow D_j^{t'}$  (resp.  $x_i^t \rightarrow x_i^{t'}$ ) has *non-zero impact* on outcome  $y^t$  iff  $\text{diff}_Y(y^t, y^{t'}) \neq \emptyset$ , where  $y^{t'}$  is the new outcome computed using  $D_j^{t'}$  (resp.  $x_i^{t'}$ ) in (1).

## 2.2. Selective re-computation steps

Using this notation, we formulate steps S1–S4 as follows. Firstly, *blind re-computation* following a change in a data dependency  $D_j$ :  $D_j^t \rightarrow D_j^{t'}$  is simply the complete re-execution of (1) for each  $y^t$  in  $Y^t$ , by replacing  $D_j^t$  with  $D_j^{t'}$ .<sup>10</sup>

**S1: Diff functions** We assume one can define a family of type-specific *data diff* functions, which quantify the extent of changes that occur over time in either  $x$ ,  $D_j$ , or  $y$ . Specifically:

$$\text{diff}_X(x_i^t, x_i^{t'}) \quad \text{diff}_Y(y_i^t, y_i^{t'}) \quad (2)$$

compute the differences between two versions of  $x_i$  of type  $X$ , and two versions of  $y_i$  of type  $Y$ . Similarly, for each source  $D_j$ ,

$$\text{diff}_{D_j}(D_j^t, D_j^{t'}) \quad (3)$$

quantifies the differences between two versions of  $D_j$ . The values computed by each of these functions are type-specific data structures, and will also depend on how changes are made available. For instance,  $D_j^t, D_j^{t'}$  may represent successive transactional updates to a relational database. More realistically in our analytics setting, and on a longer time frame, these will be two releases of  $D_j$ , which occur periodically. In both cases,  $\text{diff}_{D_j}(D_j^t, D_j^{t'})$  will contain three sets of added, removed, or updated records. The only assumption we make on these functions is that they should all report the empty difference when their inputs are identical:  $\text{diff}_T(v, v) = \emptyset$  for any type  $T$ .

**S2: Identifying the scope of change** Suppose an outcome  $y^t$  is produced a  $D_j^t$  which is later updated:  $D_j^t \rightarrow D_j^{t'}$ . Intuitively, if  $\text{exec}(P, x_i^t, D^t)$  has not used any of the data in  $\text{diff}_{D_j}(D_j^t, D_j^{t'})$ , then  $y^{t'} = \text{exec}(P, x_i^t, D^{t'}) = y^t$ , as we have assumed that  $P$  is deterministic.

Thus, we may be able to determine with certainty, required in our lossless re-computation setting, that some of the outcomes  $y^t$  do not need re-computing, provided we maintain a detailed account of exactly which data each execution of  $P$  has used, from each of its external data resources. Following this intuition, we address the problem to identify the elements from a population of outcomes  $Y^t$  that are *out of scope*, that is, those for which re-execution is certainly going to produce identical results given changes in any of its data dependencies.

**S3: Partial re-execution** Suppose  $P$  is specified as a workflow, described as a directed acyclic graph of  $k$  processing elements  $P_1 \dots P_k$ , connected through data dependencies. Given an execution of  $P$  as in (1) and changes of the form  $D_j^t \rightarrow D_j^{t'}$  and/or  $x^t \rightarrow x^{t'}$ , we want to identify the minimal subset  $P'$  of  $h \leq k$  processing elements in  $P$ , such that executing  $P'$  using  $x^{t'}$  and  $D_j^{t'}$  yields the same result as executing  $P$  entirely. Past research [14, 15] outlines conditions under which effective partial execution is viable, specifically when  $P$  has a dataflow structure.

**S4: Differential execution** Given an execution of  $P$  as in (1) and changes as above, in some cases it may be possible to refresh an outcome  $y^t$  by re-computing  $P$  using only the *differences* between old and new versions of the inputs or of the data resources. As difference sets are much smaller than the entire inputs or reference data resources, especially in the case of Big Data problems, this may result in significant savings in computation time.

## 2.3. Requirements for selective re-computation

The architectural pattern we adopt is that of a meta-process (the `ReComp` process) that can provide at least two minimal capabilities relative to an underlying process  $P$ :

- to monitor changes in the inputs, dependencies, and outputs of  $P$ , and to quantify them, e.g. by accepting type-specific  $\text{diff}_X()$ ,  $\text{diff}_D()$ , and  $\text{diff}_Y()$  functions, and
- to control the partial or entire re-execution of  $P$ .

Underpinning these capabilities are a number of technical requirements regarding collecting and storing various kinds of metadata during execution, and performing analytics on it. We specifically make use of provenance metadata, which the W3C defines as “information about entities, activities, and people involved in producing a piece of data or thing, which can be used to form assessments about its quality, reliability or trustworthiness” [19]. Specific requirements include:

- **Transparency of the process structure.** Ideally,  $P$  should be a *white box* process, that is, it should be possible to inspect its internal structure to support partial and differential execution;
- **Observability of process execution and provenance collection.** It must be possible to observe data production and consumption events that occur during execution, and use those to reconstruct the *provenance* of the outcomes;
- **Cost monitoring.** Similarly, it must be possible to assess the detailed cost (execution time, storage volume) of each execution, as this is required to learn estimates of the future cost of re-execution;
- **Process Reproducibility** Finally, it must be possible to enact a new execution of  $P$  on demand.

Not all computational models are friendly to our metadata analytics approach and satisfy all of the requirements above, either because provenance collection is available but not at the level of detail that is usable (the NoWorkflow system [20], for instance, is very good at monitoring any Python process but its provenance is too low-level to be used here), or is not provided at all. In particular, *black box* processes that do not reveal their internal structure, cost, or execution provenance, such as third party web services, are particularly hostile to our analysis.

Our experiments, however, are carried out on a platform that provides ideal support for these requirements: SVI is implemented as a workflow, providing both full transparency and coarse-grained provenance collection and cost monitoring capabilities. Furthermore, the workflow is deployed on the Azure cloud, which pro-

<sup>10</sup> Note that if the change is  $x_i^t \rightarrow x_i^{t'}$ , then trivially only the executions on input  $x_i^t$  are performed.

vides an additional mapping of resources to price, through their cost model. Exploring the extent to which our results deteriorate as the requirements above are not met is currently out of our scope.

Finally, regarding reproducibility, we note that actual re-computation of older processes  $P$  under slightly different conditions is not straightforward, as it may require redeploying  $P$  on a new infrastructure and ensuring that the system and software dependencies are maintained correctly, or that the results obtained using new versions of third party libraries remain valid. Addressing these architectural issues is a research area of growing interest [21–23], but not a completely solved problem.

### 3. Related work

To the best of our knowledge a comprehensive solution to re-computation of generic analytics processes in reaction to changes in their input data has not been discussed previously, and so the proposed meta-process is *unique*. There exist, however, a large amount of work related to particular steps which our re-computation meta-process combines. We perceive this as a chance to build our system out of existing components or at least inform our implementation of them. Note, however, that the first step, computing data difference, is rather technical and so we focus on the other three steps of our meta-process.

**The scope of re-computation.** Depending on how much insight we can have into the data analytics modules and their exact semantics, one way to partially identify scope of re-computation is to rely on fine-grained *data provenance*. Known techniques of annotating data tuples help explain why-, how- and where-provenance (see [24] for a recent survey) and can inform why the output data contain specific tuple and which input tuples contributed to produce it. While this is not enough to understand how *new* tuples would affect the output, it may still be useful to determine impact of changed and removed records. Specifically, provenance enables us to discover the subset of all input records involved in producing the output, and so a change in any of these records directly indicates impact.

Despite this being a generic and application-independent mechanism, however, two issues make it hard to use data provenance to implement impact functions. Firstly, in general we consider black-box analytics modules which not necessarily follow well understood semantics of the select, project and join operators. Secondly, our experience with genomics databases shows that changes usually involve all three types of added, removed and changed records, which limits the use of the existing data provenance techniques and makes impact analysis more difficult.

For these reasons, we have chosen to use difference functions and the coarse-grained provenance information to determine the scope of re-computation. Note that this a technical aspect is novel and not well covered in the literature.

**Partial process re-execution** has been studied extensively in the past. Perhaps the best known tool for automated re-execution of programs in reaction to any changes in their dependencies is Make.<sup>11</sup> The tool helps control the build process of a program from the program's source code. Its key feature is the ability to generate a dependency graph between source files, intermediate artifacts and outputs such that a change in one source file results in a partial rather than complete rebuild of program sources. To drive partial rebuild, Make simply uses the file modification date. Whenever any of the prerequisite files has a date newer than the target file, the relevant rule is fired off and the target file is rebuilt.

Two techniques for smart rerun, SRM, and partial process re-execution, in SPADE, are closely related to our work. Smart Rerun

Manager (SRM) [14] is part of the Kepler WFMS. The idea of *smart rerun* of a workflow, previously explored by the same group [25], is to react to changes in one or more parameters in a workflow actor by only executing those parts of the workflow that are affected by the changes, taking data dependencies into account. The approach relies on coarse-grain provenance traces and intermediate data collected and stored during workflow execution, and is derived from a similar approach implemented in VisTrails [26]. In principle, the intermediate results of workflow execution are extracted from a cache instead of being recreated by re-enacting the workflow.

Each intermediate data product is assigned a unique ID in the cache. The provenance trace is traversed from the end of the execution back to the start. For each actor found during the traversal, SRM checks whether the data products generated by this actor are still valid, i.e. are found in the cache. If that is the case, then the entire subgraph that ends with that actor does not require re-execution. In this case, the cached data is used from that point onwards.

SPADE recently implemented partial process re-execution [15]. The framework can capture fine-grained system-level provenance information and can later use it to improve effectiveness of process re-execution. By intercepting the low-level system calls, SPADE can recreate a DAG structure of the process even without explicit workflow specification. The basis of building the acyclic data dependency graph is versioning of the data artifacts. If a task within the process reads and writes to a file, every write generates a new version of the file which can potentially be reused during re-execution and rollback.

Our approach to partial re-execution is similar to both SRM and SPADE. We collect provenance of workflow-based applications like SRM does, whilst to calculate the minimal re-computation subgraph we use the data versioning mechanism provided by e-SC, which is closer to file versioning in SPADE. Also similar is that to store provenance information we use the PROV and ProvONE models, which are successors of the OPM model used by SPADE. Although the idea is not new, ours is the first implementation to operate off the e-SC workflow model, and it plays only a partial role in a more ambitious picture, where we seek to prioritise re-execution within a large collection of prior outcomes.

Techniques for **differential execution** such as incremental computation [27,28] address the problem of reacting effectively to incremental changes in the program's input data. Briefly, these techniques are based on dependency graphs, memoisation and partial evaluation – concepts similar to what we use to re-compute our process, yet applied on the scale of a single algorithm or program. A number of incremental computation solutions has also been applied to Big Data problems. Most notable are Dryadinc [29], Haloop [30] and Incoop [31] and more recently iiHadoop [32]. Again, the main difference between these and our approach is that we consider re-computation in broader sense, not limited to only a single algorithm or execution for which input data has been updated. Instead, by combining all four steps we can address the problem of selective re-computation in a comprehensive way and across many separate executions. The problem we address is to effectively reduce the number of past executions which need re-computation and also the amount of processing a single data update requires. This does not prevent the use of other incremental techniques, e.g. differential dataflows [18], iiHadoop or parallel incremental computation implemented in iThreads [33], as the basis for re-execution.

### 4. Experimental setting and blind re-computation baseline

Our experiments are based on the SVI tool, which is a natural continuation of the more complex variant calling NGS pipeline.

<sup>11</sup> <http://www.gnu.org/software/make>.

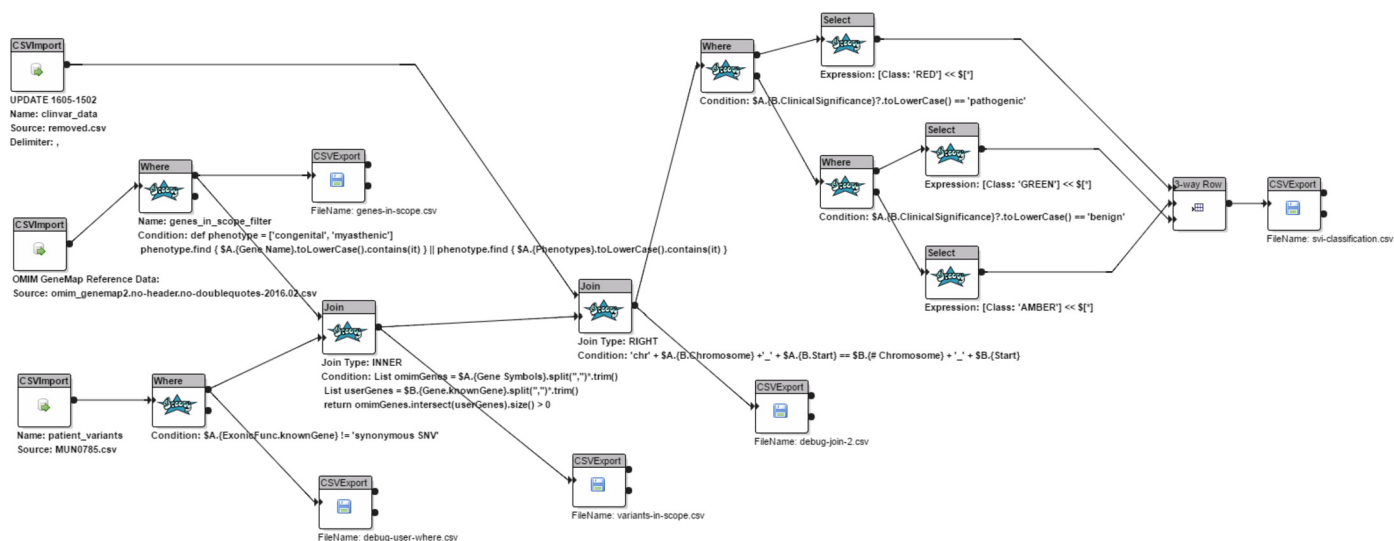


Fig. 5. The SVI tool implemented as an e-Science Central workflow.

SVI is much less resource-intensive and thus easier to work with than the complete pipeline. At the same time, it exhibits many features of the larger process: its reference databases are updated frequently and it may be run over a cohort of patients, thus it provides a very realistic example how evolving input data may influence patient’s diagnosis.

#### 4.1. The SVI workflow

SVI is implemented using the same workflow technology, namely the e-Science Central (e-SC) platform [34], a cloud-based Workflow Management System designed for scientific data management and analysis. A screenshot of the SVI implementation in e-SC is shown in Fig. 5. To recall, e-SC supports a simple dataflow programming model where processing blocks are connected to each other using data links, in a DAG topology. Our performance analysis of the variant calling workflow is described in detail elsewhere [4].

The choice of using e-SC as a platform for our experiments satisfies all of the requirements listed above. Firstly, workflows are *white box* processes, where the partial re-execution problem translates into a problem of selecting a suitable sub-graph from the whole workflow DAG structure.

Secondly, e-SC automatically records the derivation history of every workflow output from the inputs, i.e. their provenance. The provenance traces are described using the ProvONE data model [35], which extends the standard PROV data model [19]. For the purpose of human-readable description, in this paper we use the PROV-N notation [36] to present relevant provenance fragments. This includes details required for re-execution, such as the parameter settings for the processing blocks, and the version of each library and data dependency at the time of execution.

Thirdly, detailed execution costs at the level of the single processing block, as well as data storage costs, are available either through e-SC or the underlying cloud deployment (Azure, in this instance). Finally, the dependency manager that oversees the execution of e-SC workflows provides the required levels of reproducibility, i.e. the ability to re-run old workflows on demand. This rich corpus of metadata enables the kind of analysis required by our techniques, for instance estimating the cost of selecting a particular sub-graph for partial workflow re-execution.

#### 4.2. Data changes considered in the experiments

As mentioned in the introduction, for each patient SVI uses two kinds of data: the case-specific patient variant file and phenotype, and two external reference databases: OMIM GeneMap to find genes relevant to the given phenotype hypothesis, and NCBI ClinVar used to interpret the pathogenicity of patient variants.

Changes to the reference databases are very relevant because they often affect a large number of patients and are the primary cause of knowledge decay. Furthermore, the SVI reference databases change frequently, providing a good time granularity for experiments: GeneMap updates are published every day, whereas new ClinVar versions are announced every month.

In contrast, changes to patient phenotype are not considered because those represent a change, initiated by a clinician, in the actual disease hypothesis, and this automatically triggers a new investigation for the patient. Similarly, we do not consider updates to the patient variants, because those change infrequently and not enough data points would have been available in our test dataset.

#### 4.3. Experimental setup

For the purpose of this study, SVI was run on a small-scale deployment of the e-Science Central system in Microsoft Azure, consisting of the e-SC server running on a Basic A2 VM (2 CPU-cores, 3.5 GB RAM) and of a single workflow engine, hosted on a Basic A3 VM (4 CPU-cores, 7 GB RAM). Both ran Ubuntu 16.04 OS.

We used patient variants files with three phenotypes: Alzheimer’s disease, Frontotemporal dementia – Amyotrophic lateral sclerosis and CADASIL syndrome. On average they included 24 thousand records, around 39 MB in size. The OMIM GeneMap reference database was accessed on 31/Oct/2016 (unless stated otherwise) and a range of versions of NCBI ClinVar database were used, from July 2015 to October 2016. For more details about the patient variant files and reference databases please refer to Table A.6 and A.7 in Appendix A.

#### 4.4. Baseline: blind re-computation

As mentioned in the introduction, *blind re-computation* refers to the baseline case where any change at all in any of the reference databases triggers a full re-computation of the entire population of prior outcomes. Table 1 shows the effects of reacting to new versions of ClinVar regardless of the extent of the changes between



**Table 3**  
The number of records and reduction percentage of the generic and SVI-based difference sets calculated for selected versions of OMIM GeneMap. Highlighted is less favourable size reduction of the sets.

GeneMap versions $D_{old} \rightarrow D_{new}$	$ D_{new} $	$ ADDED  + 2 \cdot  CHANGED  +  REMOVED $		Reduction (%)	
		Generic $\delta$	SVI-specific $\delta$	$1 - \frac{ \delta_{gen} }{ D_{new} }$	$1 - \frac{ \delta_{SVI} }{ D_{new} }$
16-04-28 $\rightarrow$ 16-06-01	15897	27 + 196 + 1 = 224	27 + 142 + 1 = 170	98.6	98.9
16-06-01 $\rightarrow$ 16-06-02	15897	0 + 8 + 0 = 8	0 + 4 + 0 = 4	99.95	99.97
16-06-02 $\rightarrow$ 16-06-07	15910	13 + 76 + 0 = 89	13 + 52 + 0 = 65	99.4	99.6
16-06-07 $\rightarrow$ 16-10-30	16031	128 + 11944 + 7 = 12079	128 + 636 + 7 = 771	<b>24.7</b>	95.2
16-10-30 $\rightarrow$ 16-10-31	16031	0 + 10 + 0 = 10	0 + 8 + 0 = 8	99.94	99.95
16-10-31 $\rightarrow$ 16-11-01	16031	0 + 42 + 0 = 42	0 + 0 + 0 = 0	99.7	100.0
16-11-01 $\rightarrow$ 16-11-02	16031	0 + 4 + 0 = 4	0 + 0 + 0 = 0	99.98	100.0
16-11-02 $\rightarrow$ 16-11-30	16063	34 + 186 + 2 = 222	34 + 138 + 2 = 174	98.6	98.9

any two versions. The Table reports results from nearly 500 executions, concerning a cohort of 33 patients, for a total runtime of about 58.7 hours. As merely 14 relevant output changes were detected, this is about 4.2 hours of computation per change: a steep cost, considering that the actual execution time of SVI takes a little over seven minutes.

Furthermore, the table only portrays a partial picture, as it only includes reactions to monthly changes in ClinVar. A really blind approach would also react to *daily* changes to GeneMap, which are shown to have very little effect on the outcomes. For instance, comparing outputs generated using the same version of ClinVar and four consecutive versions of GeneMap: 16-10-30, 16-10-31, 16-11-01 and 16-11-02 shows that none of the patient variants was affected at all by these small changes.

The sparsity of the table should come as no surprise, as changes in the reference databases are dispersed across the whole human genome and so the chance that they may affect a particular patient are relatively small. Further details on these experiments can be found in supplementary material sheet CV-blind.

In the next sections we discuss in detail the re-computation meta-process applied to our model analytics process SVI. We start from computing data differences (S1) and then grow the amount of processing going backwards from experimenting with differential execution (S4), to applying partial re-execution (S3), up to identification of the scope of change (S2). As a result, our final experiment includes a realisation of the re-computation meta-process for the SVI tool with steps S1–S3 applied.

## 5. Data differences

The reference databases that SVI uses are in the “well-behaved” category of simple relational tables, making it easy to express differences in terms of set operations. Specifically, the added and removed subsets are just set difference between two versions, while the changed subsets are an intersection followed by a selection. The following SQL-like pseudocode specifies these operations more formally on two versions  $D_1, D_2$  of a data table:

```
ADDEDD1→2 = select * from D2
           where not exists (
             select 1 from D1
             where D1.KEY = D2.KEY
           )
```

```
REMOVEDD1→2 = select * from D1
            where not exists (
              select 1 from D2
              where D1.KEY = D2.KEY
            )
```

```
CHANGEDD1→2 = select D2.* from
              D1 inner join D2 on D1.KEY = D2.KEY
              where D1.NON-KEYc1 <> D2.NON-KEYc1 or
                 D1.NON-KEYc2 <> D2.NON-KEYc2 or ...
```

These operators assume that we have selected key attributes (possibly compound) for  $D$ . Also,  $CHANGED_{D1 \rightarrow 2}$  denotes the new version of the updated records, whereas analogous  $CHANGED_{D2 \rightarrow 1}$  is used to compute the old version of the updated records. In effect:

$$\delta^+ = ADDED_{D1 \rightarrow 2} \cup CHANGED_{D1 \rightarrow 2}$$

$$\delta^- = REMOVED_{D1 \rightarrow 2} \cup CHANGED_{D2 \rightarrow 1}$$

Note that the CHANGED operator captures all changes in any of the non-key attributes for each record. While this is generic, it ignores the meaning of the attributes relative to the specific processing and is likely to result in a large number of changes irrelevant to the process. For example, two GeneMap records that differ only in the `Comments` attribute would be flagged as different, although the comments are not used anywhere in SVI. Similarly, the only changes in ClinVar records that are relevant to SVI are those in the `ClinicalSignificance` attribute, which drive the classification of variants in the SVI output.

Thus, with the knowledge of the specific use of a relational dataset that the process makes, we partition the attributes into the `KEY`, `USED`, and `UNUSED` datasets. The CHANGED operator can then be rewritten as:

```
CHANGEDD1→2 = select D2.* from
              D1 inner join D2 on D1.KEY = D2.KEY
              where D1.U1 <> D2.U1 or
                 D1.U2 <> D2.U2 or ...
```

where  $U_i \in \text{USED}$ .

There is an obvious benefit in efficiency resulting from this more aggressive filtering of the difference sets, as illustrated in Tables 3 and 4. The tables report on the number of records of the complete GeneMap and ClinVar datasets and the difference sets calculated using the generic and SVI-specific diff operators. Using the SVI-specific operators the reduction in size is almost always about 90% or over. The only exceptions are the differences between version Jul  $\rightarrow$  Aug 2015 of ClinVar which faced a significant change at the time. Then, the SVI-specific operator yielded a reduction of 49.6%.

More limited gain is achieved when using the generic *diff* operators. In three cases the total size of the difference sets was *larger* than the new version of the ClinVar database. Similarly, the differences between GeneMap 16-06-07 and 16-10-30 computed by the

**Table 4**

The number of records and reduction percentage of the generic and SVI-based difference sets calculated for 16 versions of ClinVar. Highlighted are less favourable size reductions of the sets.

ClinVar versions $D_{old} \rightarrow D_{new}$	$ D_{new} $	$ ADDED  + 2 \cdot  CHANGED  +  REMOVED $		Reduction (%)	
		Generic $\delta$	SVI-specific $\delta$	$1 - \frac{ \delta_{gen} }{ D_{new} }$	$1 - \frac{ \delta_{SVI} }{ D_{new} }$
15-07 → 15-08	252656	35087 + 425794 + 85987 = 546868	35087 + 6302 + 85987 = 127376	<b>-116.4</b>	<b>49.6</b>
15-08 → 15-09	259714	7273 + 16952 + 215 = 24440	7273 + 1342 + 215 = 8830	90.6	96.6
15-09 → 15-10	262498	2832 + 11888 + 53 = 14773	2832 + 1174 + 53 = 4059	94.4	98.5
15-10 → 15-11	277902	15550 + 108588 + 146 = 124284	15550 + 4300 + 146 = 19996	55.3	92.8
15-11 → 15-12	279174	1376 + 489530 + 104 = 491010	1376 + 472 + 104 = 1952	<b>-75.9</b>	99.3
15-12 → 16-01	280379	1523 + 23740 + 318 = 25581	1523 + 2224 + 318 = 4065	90.9	98.6
16-01 → 16-02	285041	4710 + 26304 + 48 = 31062	4710 + 1490 + 48 = 6248	89.1	97.8
16-02 → 16-03	286684	2477 + 235330 + 453 = 238260	2477 + 2510 + 453 = 5440	<b>16.9</b>	98.1
16-03 → 16-04	290432	3855 + 27088 + 107 = 31050	3855 + 1282 + 107 = 5244	89.3	98.2
16-04 → 16-05	290815	858 + 15732 + 475 = 17065	858 + 1158 + 475 = 2491	94.1	99.1
16-05 → 16-06	306503	18004 + 81738 + 2298 = 102040	18004 + 7174 + 2298 = 27476	66.7	91.0
16-06 → 16-07	320469	14496 + 56692 + 530 = 71718	14496 + 6696 + 530 = 21722	77.6	93.2
16-07 → 16-08	326856	6558 + 58238 + 174 = 64970	6558 + 31356 + 174 = 38088	80.1	88.3
16-08 → 16-09	327632	1020 + 18838 + 244 = 20102	1020 + 1104 + 244 = 2368	93.9	99.3
16-09 → 16-10	349074	22758 + 654486 + 630 = 677874	22758 + 13228 + 630 = 36616	<b>-94.2</b>	89.5

generic operators were only 24.7% smaller than the new version of the database. Clearly, in such cases it is more effective to ignore the difference sets and use only the new version of the data.

Another important aspect of calculating the difference sets is the changing set of attributes. For example, the ClinVar attributes have changed three times since February 2015. These changes in the schema disrupt our difference operators because the three sets KEY, USED, UNUSED change, and also they are no longer perfectly aligned across versions. Therefore, in our implementation of ClinVar diff we assumed that we would compare only columns common in both versions and ignore the added and removed columns. Currently, our SVI re-computation supports any version of ClinVar since Feb 2015.

## 6. Differential execution

Given the difference sets of the changed inputs we can look at re-execution of  $P$  using merely  $\text{diff}_D(d^t, d^{t'})$  – the differences between two versions of (one or more) reference dataset  $D$ . Some of these ideas are grounded in prior research on the incremental computation and differential computation domains [37,18].

Using SVI as our testbed, we show that under some conditions this is feasible to do and results in substantial savings. However, in the general case  $P$  requires modifications in order to yield a valid result.

### 6.1. Computing on data versions differences

To make the idea precise consider our baseline execution (1):

$$\langle y^t, c^t \rangle = \text{exec}(P, x^t, D^t) \quad (4)$$

We are now going to focus on changes to  $D$ , thus we assume  $x^t$  is constant over time:  $x^{t'} = x^t = x$  (in SVI, this means we consider one patient at a time). For simplicity of exposition, initially we also assume a single  $D$  with states  $D^t, D^{t'}$ . In the common case where  $D$  is a relation and  $D^t$  consists of a set of records, such as a CSV-formatted file (the case for ClinVar), we can express  $\text{diff}_D(D^t, D^{t'})$  in terms of set differences:

$$\text{diff}_D(D^t, D^{t'}) = \langle \delta^+, \delta^- \rangle$$

where:

$$D^{t'} = D^t \setminus \delta^- \cup \delta^+ \quad (5)$$

and  $\delta^+$  denotes the added records and *new version* of updated records whilst  $\delta^-$  are records removed from  $D^t$  and the old version of the records that are going to be updated. Note that in the case of ClinVar and most bioinformatics databases,  $\delta^-$  include *retractions*, which are much less frequent than additions of new records.

Our contention is that the computation of new outcome

$$\langle y^{t'}, c^{t'} \rangle = \text{exec}(P, x, D^{t'}) \quad (6)$$

can be broken down into two smaller computations that only use  $\delta^+, \delta^-$  and produce partial outcomes  $y_{+}^{t'}, y_{-}^{t'}$ , which can then be combined with  $y^t$  to yield  $y^{t'}$ . This may require an additional  $\text{merge}(\cdot)$  function that is process-specific. More precisely, we break (6) down into:

$$\langle y_{+}^{t'}, c_{+}^{t'} \rangle = \text{exec}(P, x, \delta^+) \quad (7)$$

$$\langle y_{-}^{t'}, c_{-}^{t'} \rangle = \text{exec}(P, x, \delta^-) \quad (8)$$

$$\langle y^{t'}, c_m^{t'} \rangle = \text{merge}(y^t, y_{+}^{t'}, y_{-}^{t'}) \quad (9)$$

This breakdown is beneficial if the resulting total cost is less than  $c^{t'}$ :

$$c_{+}^{t'} + c_{-}^{t'} + c_m^{t'} < c^{t'}$$

Firstly, let us consider the case in which  $P$  implements a “well-behaved” function that is distributive over set union and difference. Using (5) and (6) and ignoring cost for the time being, we can write:

$$\begin{aligned} y^{t'} &= \text{exec}(P, x, D^{t'}) \\ &= \text{exec}(P, x, D^t \setminus \delta^- \cup \delta^+) \\ &= \text{exec}(P, x, D^t) \setminus \text{exec}(P, x, \delta^-) \cup \text{exec}(P, x, \delta^+) \\ &= y^t \setminus y_{-}^{t'} \cup y_{+}^{t'} \end{aligned} \quad (10)$$

Thus, in this case,  $y_{+}^{t'}, y_{-}^{t'}$  can be automatically combined into  $y^{t'}$ . However, distributivity is a strong assumption, which does

not hold for many practical cases. For example, SVI as a whole distributes only over set difference and union of selected inputs. But in the general case,  $P$  may need to be modified in order to combine the partial results using an ad hoc *merge* function.

To illustrate this situation note that SVI essentially consists of four steps (cf. Fig. 2):

- a) SELECTION operation that selects from GeneMap only genes relevant to user-defined phenotype  $ph$ , producing genes in scope  $GS$ :

$$GS = \sigma_{ph}(\text{GeneMap})$$

- b) INNER JOIN operation between the input variants,  $x$ , and the result of the previous query  $GS$ , producing variants in scope  $VS$ :

$$VS = x \bowtie GS$$

- c) RIGHT OUTER JOIN to combine pathogenicity annotations from ClinVar with corresponding  $VS$  yielding  $VS_p$ :

$$VS_p = \text{ClinVar} \bowtie VS$$

- d) final classification into the traffic light system, which adds new classification column to  $VS_p$ :

$$\text{classify}(VS_p)$$

Given these steps, SVI is specified by the following expression:

$$\begin{aligned} \text{SVI}(x, ph, \text{GeneMap}, \text{ClinVar}) \\ = \text{classify}(\text{ClinVar} \bowtie (x \bowtie \sigma_{ph}(\text{GeneMap}))) \end{aligned}$$

Note, however, that (a), (b) and (d) are all distributive over set union and difference, whereas (c), the right outer join, distributes only for the right-hand side argument. Therefore, we can automatically combine partial outcomes when  $\delta^+$  and  $\delta^-$  are computed for GeneMap,  $GS$ ,  $x$ ,  $VS$ , or  $VS_p$ , whilst differences in ClinVar require a custom merge function.

Regarding GeneMap, the computation steps are as follows. Let

$$y^t = \text{SVI}(x, ph, GM^t, CV)$$

be the original computation, and

$$GM^t = GM^t \setminus \delta^- \cup \delta^+$$

be a new version of GeneMap, expressed in terms of version differences. The new  $y^t$  can be computed as:

$$\begin{aligned} y^t &= \text{SVI}(x, ph, GM^t, CV) \\ &= \text{classify}(CV \bowtie (x \bowtie \sigma_{ph}(GM^t))) \\ &= \text{classify}(CV \bowtie (x \bowtie \sigma_{ph}(GM^t \setminus \delta^- \cup \delta^+))) \\ &= \text{classify}(CV \bowtie (x \bowtie (\sigma_{ph}(GM^t) \setminus \sigma_{ph}(\delta^-) \\ &\quad \cup \sigma_{ph}(\delta^+)))) \\ &= \text{classify}(CV \bowtie ((x \bowtie \sigma_{ph}(GM^t)) \\ &\quad \setminus (x \bowtie \sigma_{ph}(\delta^-)) \cup (x \bowtie \sigma_{ph}(\delta^+)))) \quad (11) \\ &= \text{classify}(CV \bowtie (x \bowtie \sigma_{ph}(GM^t))) \\ &\quad \setminus \text{classify}(CV \bowtie (x \bowtie \sigma_{ph}(\delta^-))) \\ &\quad \cup \text{classify}(CV \bowtie (x \bowtie \sigma_{ph}(\delta^+))) \\ &= \text{SVI}(x, ph, GM^t, CV) \\ &\quad \setminus \text{SVI}(x, ph, \delta^-, CV) \cup \text{SVI}(x, ph, \delta^+, CV) \\ &= y^t \setminus y_-^t \cup y_+^t \end{aligned}$$

Although the same approach does not work for changes in ClinVar, we can still adapt SVI and define a bespoke *merge*() function to combine partial results in a way that is *semantically meaningful* for the specific data and process. An implementation of such a function would filter the results from the right outer join  $y_-^t$  and  $y_+^t$  by removing rows with null in ClinVar columns, essentially turning right outer join into inner join:

$$z_-^t = \text{filter}(y_-^t) \quad z_+^t = \text{filter}(y_+^t) \quad (12)$$

Then, using the filtered products it performs specialised set operations:

$$y^t = y^t \cup_{\text{amb}} z_-^t \cup_{\text{wrt}} z_+^t \quad (13)$$

where  $a \cup_{\text{amb}} b$  sets the classification to amber for all rows in  $a$  that match  $b$ , whereas  $a \cup_{\text{wrt}} b$  overwrites the classification for all rows in  $a$  that match  $b$  on non-ClinVar columns; both operations leave non-matching rows intact. Given definitions (12) and (13), we can define the specialised merge function as:

$$\begin{aligned} \text{merge}_{\text{SVI}}(y^t, y_+^t, y_-^t) \\ = y^t \cup_{\text{amb}} \text{filter}(y_-^t) \cup_{\text{wrt}} \text{filter}(y_+^t) \end{aligned} \quad (14)$$

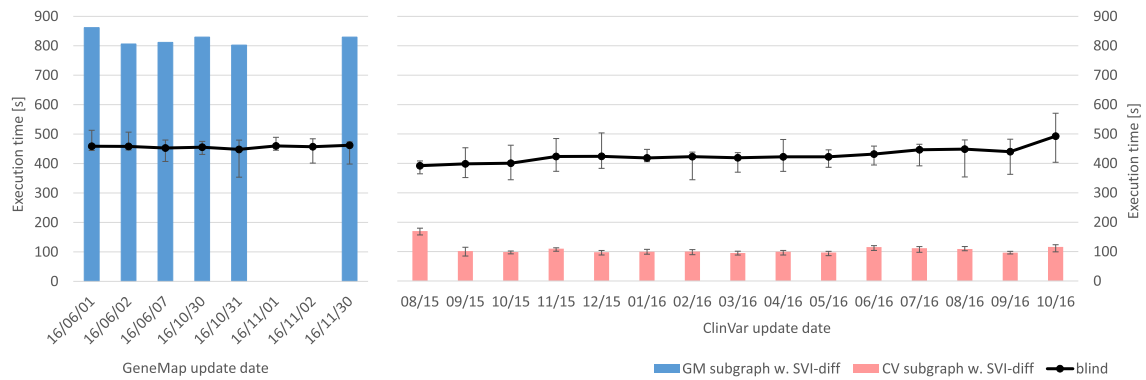
As we can see, this approach may need a substantial amount of process refactoring and makes the definition of *merge* encode some of the process semantics to operate on data differences. Nevertheless, in the rest of this section we are going to illustrate the practical steps in computing the differences  $\delta^+$ ,  $\delta^-$  and the partial outputs  $y_+^t$  and  $y_-^t$  on SVI, which will be useful to address the scope analysis.

## 6.2. Re-computation using the difference sets

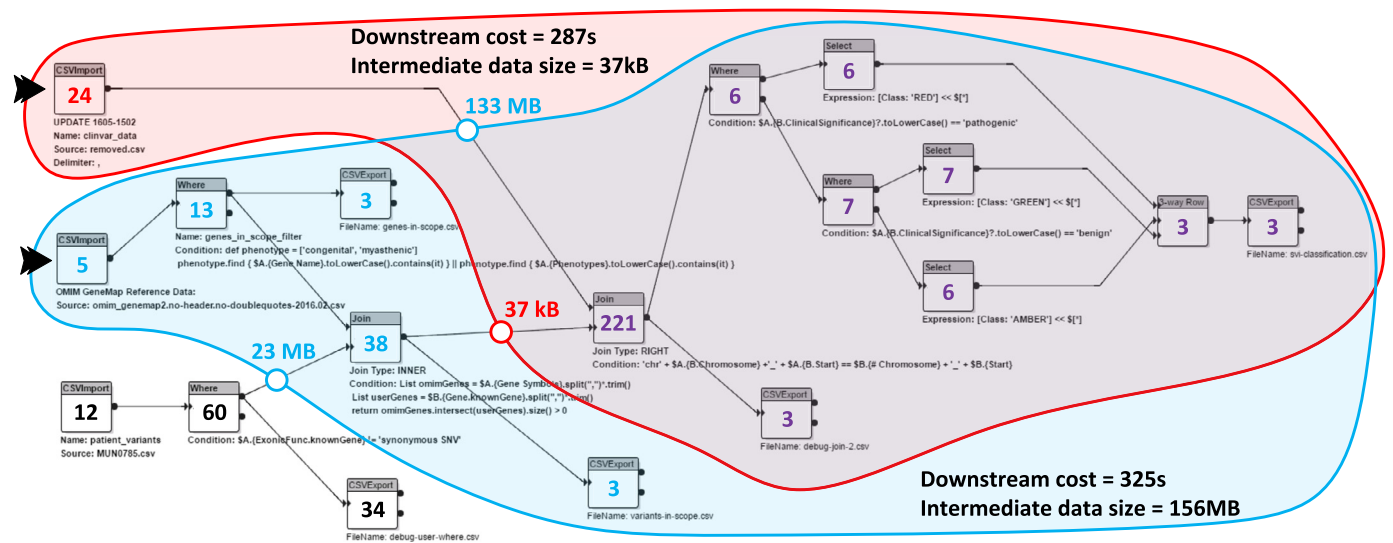
To see the effect of using the difference sets on runtime we executed SVI with the range of GeneMap and ClinVar difference sets shown in Table 3 and 4. Note that in the case of ClinVar differences we did not include the merge function defined earlier in (14). However, doing so would not affect the runtime significantly as the SVI outputs contain only about a dozen rows. Fig. 6 shows the execution times.

Interestingly, the results indicate two very distinct cases. First, using the difference sets to calculate the output yields clear runtime savings for changes in ClinVar. Re-computation time oscillated around 100 seconds with the only exception for the considerable changes between the July and August 2015 versions of the database (cf. Table 4). However, using the GeneMap difference sets we observed loss in the execution time in most cases. Even if the changes were minimal (e.g. 16-06-01→16-06-02 and 16-06-02→16-06-07) and the difference sets contained only a few records, re-execution took about 400 seconds for  $\delta^+$  and  $\delta^-$  separately; over 800 seconds altogether. In two cases we could skip re-execution because the difference sets were empty.

That problem with GeneMap differences stems from the fact that this database is nearly two orders of magnitude smaller than ClinVar. Therefore, the majority of runtime is spent on tasks processing ClinVar whilst the smaller GeneMap file does not affect overall execution time that much. We observed some savings only for two blocks: the WHERE and JOIN located at the front of the pipeline. But the remainder of the pipeline used the complete ClinVar database. Conversely, when ClinVar undergoes changes, the data is used by SVI at the tail of the pipeline where the longest running JOIN block is located (cf. Fig. 8). Thus, using the difference sets rather than the complete version of ClinVar, we could lower the runtime of that block significantly and reduce the total execution time of the relevant workflow subgraph.



**Fig. 6.** The re-computation time of the SVI workflow using the difference sets following changes in GeneMap (left; ClinVar version 16-08) and ClinVar (right; GeneMap version 16-10-31).



**Fig. 7.** The illustration how retrospective provenance trace of a complete execution can be used to calculate the cost of partial execution; the black arrow at the top-left corner indicates the starting block; the red circle indicates the required intermediate data; red numbers in the workflow blocks denote execution time in seconds. (For interpretation of the colours, the reader is referred to the web version of this article.)

Overall, even if the use of difference sets can reduce runtime of a single execution to some extent, the savings depend on the structure of the process and may not be enough to compensate for the fact that two executions are needed – one for  $\delta^+$  and one for the  $\delta^-$  difference set.

**7. Partial re-execution**

The third step in our re-computation model is partial re-execution. As shown earlier in Fig. 4, to implement it we do not require the actual difference data but only information about the data dependencies. In Sec. 4.1 we mentioned that e-SC generates one ProvONE-compliant provenance trace for each workflow run. We exploit these traces to identify the minimal sub-workflow that is affected by the change [14,15].

Suppose we record a change of the form  $d^t \rightarrow d^{t'}$  in reference data, and let  $I$  be a past invocation of our workflow. The source blocks for any sub-workflow that is affected by the change are those activities  $A$  that were executed as part of  $I$  and that used  $d^t$  directly. These can be obtained from the query:

`:- wasPartOf(A, I), used(A, dt)`

Note that the change event itself can be recorded using a provenance assertion:

`wasDerivedFrom(dt', dt)`

In this case, the query becomes:

`:- wasDerivedFrom(dt', D), wasPartOf(A, I), used(A, D)`

where  $D$  is now a variable that represents the previous version of  $d^t$ .

Having determined the source blocks, we expand the workflow recursively, by traversing the provenance graph for invocation  $I$ , downstream. At each step we seek two possible patterns:

1. `execution(A1), execution(A2), wasInformedBy(A2, A1)`: given  $A_1$ , find all activities  $A_2$  that have been triggered by  $A_1$ . This pattern represents a connection from  $A_1$  to  $A_2$ , where the intermediate data that flows over the link during execution is implicit;
2. `execution(A1), execution(A2), wasGeneratedBy(D, A1), used(A2, D)`: Here the dependency between  $A_1$  and  $A_2$  is represented explicitly by the intermediate data product,  $D$ .

Fig. 7 shows the sub-workflows related to a change in GeneMap (blue area) and ClinVar (red area). The black arrows on the left indicate the starting blocks for the sub-workflows. The overlapping area between the two sub-workflows contains the blocks that are affected by either of the two changes. Clearly, a partial execution following a change in only one of the databases requires that the

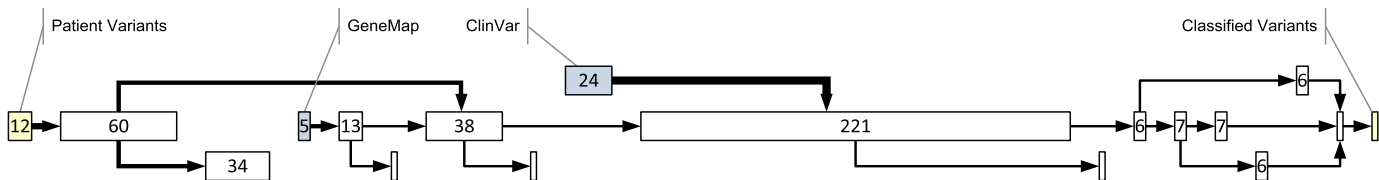


Fig. 8. One of the possible schedules of the SVI workflow tasks; wider arrows denote more data flowing between tasks; numbers in boxes represent task execution time in seconds (ClinVar  $v = 16-09$ , GeneMap  $v = 16-10-31$ , PV = B\_0201).

intermediate data at the boundary of the blue and red areas be cached.

To provide a measure of the trade-off between additional space requirements and the savings in execution time, we have annotated the figure with the execution time (in seconds) for each of the blocks, and with the size of the cached intermediate data. For example, for a GeneMap change, the corresponding partial execution took about 325 seconds (or 5 m:25 s), whereas a ClinVar change required 287 seconds (or 4 m:47 s) to re-execute. Recalling that the cost of a complete execution was 455 seconds (or 7 m:35 s), these are savings of 28.5% and 37%, respectively. The corresponding additional storage costs are 156 MB and just 37 KB, resp. The complete results of the partial re-computation following changes in the reference databases are included in supplementary material sheets CV-subgraph and GM-subgraph.

We note that these figures are obtained from the provenance traces, namely using the standard `prov:startTime` and `prov:endTime` properties of activity, along with an additional `recomp:dataSize` property, which we added to record the size of the entities transferred between blocks.

Fig. 8 provides an alternative view of a possible schedule for the same workflow, which shows the parts of the workflow affected by each of the changes, along with the rendering of the execution times and amount of data involved. We can see that a change in ClinVar affects only a sub-workflow starting in the middle of the workflow, whereas a change in GeneMap affects almost all of the blocks. However, both sub-workflows include the longest running block, which limits the amount of savings that can be achieved.

## 8. Identifying the scope of change

We now can address the second step (S2) from Sec. 2, namely how to identify the scope of a change in reference data  $D$  that is used to produce a large population,  $Y$ , of outcomes [3]. As mentioned, SVI is once again a good case study for this problem as the same process is executed over a possibly large cohort of patients (thousands). Whilst these executions are all independent of one another, they all depend on the same reference datasets. The *scope* of a change in any of these dependencies  $D$  is subset  $Y_s \subseteq Y$  of outcomes affected by change  $D^t \rightarrow D^{t'}$ .

Provided that  $\text{diff}_D(D^t, D^{t'})$  includes only changed and removed records and the process consist of a set of known transformations, fine-grained provenance solutions such as [38–40] can help identify whether or not  $y \in Y_s$ . Briefly, fine-grained data provenance enables both forward and backward queries. A forward query  $\phi(i)$  retrieves the output records that are associated by derivation to an input record  $i$ , while a backwards query traces the input records that contributed to a given output. Thus, given  $\delta^-$  and assuming that  $\text{ADDED}_{D^t \rightarrow D^{t'}} = \emptyset$ , the problem of finding whether  $y \in Y_s$  can be formulated using the forward query  $\phi$  as follows:

$$y \in Y_s \iff \exists i \in \delta^- : \phi(i) \cap y \neq \emptyset \quad (15)$$

Effectively, the forward query  $\phi$  plays the role of an indexing mechanism that maps input records to the outputs they contributed to. There are two main limitations of this approach, how-

ever. Firstly, the fine-grained data provenance is unable to handle newly added records, because there is simply no provenance trace about these records until the process runs and produces the output. Secondly, while fine-grained provenance can capture a rich set of transformations used in practice (see, e.g. [41,40]), it will not be available for black-box analytics processes that do not reveal their internal sub-processes. In this case, provenance will naively report that all input records contributed to produce each output record, with a trivial scope that includes all outputs.

A possible statistical approach to establishing whether  $y \in Y_s$  with some confidence is to sample a number of prior  $y$  from  $Y$ , compute the corresponding  $y'$ , and use the differences  $\text{diff}_Y(y, y')$  to try and learn an estimator for the differences on the unobserved new outcomes. This approach, however, is likely to be sensitive to the specific types of data and process involved and may not always yield robust estimators.

### 8.1. The basic scoping algorithm

Instead, we propose a scope determination algorithm that relies on the coarse-grained provenance associated with past runs to determine which outcomes  $y$  have used a version of  $D$ . Coarse-grained provenance, however, only indicates whether or not a dependency on  $D$  existed, but not which specific data from version  $D^t$  of  $D$  was used. It is, therefore, possible that an outcome  $y$  that depended on  $D$  is not really in the scope of change  $D^t \rightarrow D^{t'}$ , for instance because the process used data from  $D^t$  that has not changed in  $D^{t'}$ . These are *candidate invocations* which must be further analysed to determine the actual impact of change on each of them.

To carry out this analysis further, we propose to re-execute  $P$  one task at a time using the difference sets ( $\delta^-$  and  $\delta^+$ ). We first consider the case in which the tasks of  $P$  are distributive over set union and difference, as described in (10) and (11). Then, we can execute the tasks one at a time until either we observe an empty result or we reach the end of the process. In the former case we know that the invocation was out of scope and no full re-execution is needed. In the latter case we can combine the original output with the final  $\delta^-$  and  $\delta^+$  to obtain the updated result of  $P(D^{t'})$ . Clearly, this approach is beneficial if computing  $P$  on the difference sets is faster than computing  $P$  using the entire  $D^{t'}$ . Algorithm 1 formalises this approach.

Procedure `SELECTIVEEXEC` takes process  $P$  and two versions of its input data together with the coarse-grained provenance information represented by *history database*  $H$ . First, difference sets  $\delta^+$  and  $\delta^-$  are computed between the two versions of the input data (line 2).  $H$  is then queried in line 3 to find occurrences of statements of form:

`used( $a, D^t$ ), wasPartOf( $a, I$ ), wasAssociatedWith( $I, \dots, P$ )`

indicating that  $D^t$  was used by a specific activity  $a$  that was part of execution  $I$  of process  $P$ , or

`used( $I, D^t$ ), wasAssociatedWith( $I, \dots, P$ )`

**Algorithm 1** Simple selective re-computation of a population of invocations of the process that is distributive over set union and difference.

```

1: procedure SELECTIVEEXEC( $P, D^t, D^t, H$ )
2:    $\langle \delta^+, \delta^- \rangle = \text{diff}_D(D^t, D^t)$ 
3:    $\mathcal{I} \leftarrow \text{LISTINVOCATIONS}(H, P, D^t)$ 
4:   for all  $I \in \mathcal{I}$  do
5:      $G \leftarrow \text{MINIMALSUBGRAPH}(I, D^t)$ 
6:     while  $G \neq \emptyset$  and  $(\delta^+ \cup \delta^-) \neq \emptyset$  do
7:        $\text{task} \leftarrow \text{POP}(G)$ 
8:        $\delta^+ \leftarrow \text{task}(\delta^+)$ 
9:        $\delta^- \leftarrow \text{task}(\delta^-)$ 
10:    end while
11:    if  $G = \emptyset$  then
12:       $y \leftarrow \text{GETOUTPUT}(\text{task})$ 
13:       $\text{SETOUTPUT}(P, D^t) \leftarrow y \setminus \delta^- \cup \delta^+$ 
14:    else
15:       $y \leftarrow \text{GETOUTPUT}(P, D^t)$ 
16:       $\text{SETOUTPUT}(P, D^t) \leftarrow y$ 
17:    end if
18:  end for
19: end procedure

```

indicating, more broadly, that  $D^t$  was used at some unspecified point during  $I$ . In both cases the provenance traces identify the set of candidate invocations.

Each of these candidate invocations is then re-executed using the difference sets (loop in lines 4–14). To do it efficiently the algorithm computes the *minimal subgraph* of  $I$  that needs re-computation (line 5), as discussed earlier in Sect. 7; for SVI it is one of the subgraphs shown in Fig. 7. The inner loop in lines 6–10 walks through each task of the minimal downstream graph (in topological order) and re-executes the task using the difference sets until either both partial outputs are empty or all tasks have been visited. Assuming that the tasks are distributive under set union and difference, the latter case allows us to generate output of  $P(D^t)$  using the previous output and partial outcomes of the last task of  $P$  (lines 12–13).

For simplicity of presentation the presented algorithm can only work for a linear graph of tasks. Nonetheless, making it work for an arbitrary directed acyclic graph with multiple entry points for  $D^t$  is a straightforward extension.

## 8.2. Practical realisation of scoping

The main limitation of Algorithm 1 is that it may only be applied across tasks that distribute over set union and difference. That is a strong assumption which is challenging even for a simple example like SVI. Likewise, it is challenging in the much more complex case of NGS pipelines in which the alignment tools (e.g. `bwa`, `samtools`) need access to the complete human reference genome. They cannot perform sequence realignment if given only a difference between two versions of the reference genome. Thus, to benefit from this algorithm in practice we extended it to allow more diverse process tasks.

Algorithm 2 presents the inner while loop which includes a set of additional checks to make sure that only tasks which can use difference sets properly are considered. Lines 9–11 handle the case from Algorithm 1 – distributive tasks. Then, in lines 12–15, the algorithm tries to use the incrementalised version of a task if one is available. That might handle the non-distributive right outer join in SVI following an approach proposed e.g. by [42]. However, as implementing an incremental version of a task is known to be a difficult problem in general, our algorithm includes one more case which we explore in more detail below.

**Algorithm 2** An extension of the selective re-computation over the executions dimension to handle various types of computing tasks. The while loop in Algorithm 1 may be changed as follows.

```

6: ...
7: while  $G \neq \emptyset$  and  $(\delta^+ \cup \delta^-) \neq \emptyset$  do
8:    $\text{task} \leftarrow \text{POP}(G)$ 
9:   if  $\text{ISDISTRIBUTIVE}(\text{task})$  then
10:     $\delta^+ \leftarrow \text{task}(\delta^+)$ 
11:     $\delta^- \leftarrow \text{task}(\delta^-)$ 
12:   else if  $\text{ISINCREMENTALIZED}(\text{task})$  then
13:     $\text{incTask} \leftarrow \text{GETINCREMENTAL}(\text{task})$ 
14:     $\delta^+ \leftarrow \text{incTask}(\delta^+)$ 
15:     $\delta^- \leftarrow \text{incTask}(\delta^-)$ 
16:   else
17:     $\text{impTask} \leftarrow \text{GETIMPACTFUNCTION}(\text{task})$ 
18:     $y \leftarrow \text{GETOUTPUT}(\text{task})$ 
19:    if  $\text{impTask}(\delta^+, \delta^-, y) = \text{true}$  then
20:       $\text{PUSH}(G, \text{task})$ 
21:       $\text{tmp}_d \leftarrow \text{GETINPUT}(\text{task})$ 
22:       $\text{EXECUTEWORKFLOW}(G, \text{tmp}_d)$ 
23:      return
24:    else
25:       $\text{PUSH}(G, \text{task})$ 
26:      break
27:    end if
28:   end if
29: end while
30: ...

```

Lines 17–27 handle all other tasks for which we first obtain an impact function. The impact function cannot compute the actual output of the task given its partial input  $\delta^-$  and  $\delta^+$ . It can, however, determine whether the partial input is *likely* to affect output  $y$  of the task. Briefly, the *impT* function returns true to denote that the partial input has non-zero impact on the output, and false otherwise. Given that, if the function returns true, the algorithm returns the current task back to the front of  $G$  and re-executes the subworkflow using the complete past input of the task (lines 20–23). Afterwards the algorithm is completed and we can return from the procedure.

In the case there is no impact of partial input on the task output we make sure  $G$  is non-empty and break the inner loop (lines 25–26). Consequently, the procedure can set the output of  $P(D^t)$  using the previous workflow output (Algorithm 1, lines 14–16).

Although for all tasks that cannot work with partial inputs the proposed algorithm forces us to use an impact function, in the simplest *default* implementation it may always return true to indicate that any change in the input affects the output. That is likely to result in more re-computation than needed but allows any arbitrary task to be handled correctly while giving chance to provide impact functions whenever possible. For example for SVI and the problematic right outer join we were able to use *inner join* as an accurate impact function.

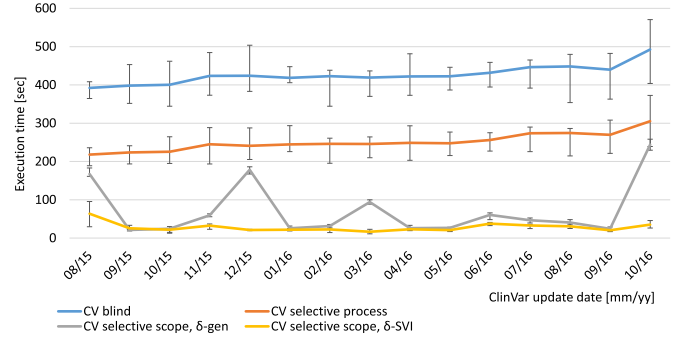
## 8.3. Scoping effectiveness

Regarding the effectiveness of this algorithm, note that the extended version of our algorithm can reduce the amount of work by making the following assumptions. First, the use of difference sets to calculate task output (lines 10–11 and 14–15) is much faster than when using the complete input data. Second, the output of these re-executions is likely to return an empty response, and so the inner while loop can terminate early with  $G \neq \emptyset$ . Third, the number of non-distributive and non-incrementalized tasks is small or, alternatively, the provided impact functions are fast, accurate and more effective than the default ‘return true’ implementation.

**Table 5**

Changes observed in the output of the SVI tool when executed with the difference sets computed for NCBI ClinVar reference database using the SVI specific  $\delta$  function; ■ denotes the need for re-execution with the complete new version of ClinVar ( $D_{ac} \neq \emptyset$  or  $D_r \neq \emptyset$ ), ‘.’ denotes only task re-execution with the difference sets ( $D_{ac} = \emptyset$  and  $D_r = \emptyset$ ).

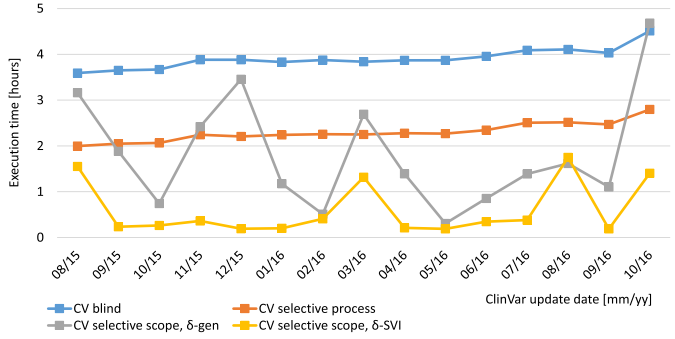
Phenotype hypothesis	Frontotemporal Dementia- Amyotrophic Lateral Sclerosis										Alzheimer's disease																											
	Variant file	D_1071	D_1049	D_1041	D_0899	D_0854	D_0830	C_0171	C_0098	C_0056	C_0053	C_0051	B_0307	D_1136	C_1457	C_0072	C_0071	C_0068	C_0065	B_0396	B_0384	B_0370	B_0365	B_0358	B_0338	B_0331	B_0229	B_0214	B_0209	B_0208	B_0203	B_0202	B_0201	B_0198				
08/15	■	.	.	.	.	.	.	.	.	.	.	■	.	.	.	■	.	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■		
09/15	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
10/15	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
11/15	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
12/15	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
01/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
02/16	.	.	■	.	■	.	.	.	.	.	.	■	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	
03/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
04/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
05/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
06/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
07/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
08/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	
09/16	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.	.
10/16	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	



**Fig. 9.** Minimum, average and maximum execution time of the SVI workflow per patient variant file in four approaches to re-computation.

Noting that all these assumptions are valid for SVI, we tested the hypothesis that the approach is indeed beneficial. We show in Table 5 that running the process using the proposed algorithm and the SVI-specific diff function we were able to avoid the majority of re-computations which used the complete new ClinVar version. We reduced the number of complete re-executions of the workflow from 495 down to 71. In Table A.8 in the appendix we show also the re-computation matrix for the algorithm which used the generic diff function. In that case the reduction was less significant and required 302 complete re-executions. That is because the generic diff searched for changes in every single column of the ClinVar data, most of which were irrelevant to SVI.

Figs. 9 and 10 show the effect of running our algorithm on the re-computation time. The former presents the average time required to re-compute a single patient variant file. For the majority of cases running SVI with the difference sets was much quicker than with complete ClinVar data. In a few cases, e.g. when using the difference between the versions from September and October 2016, some re-executions were slightly slower than the partial re-computation. That was due to extensive changes in the ClinVar database at the time and so almost all rows were reported as changed. This did not occur, however, when using the SVI-specific diff function. Then, the total time was significantly lower than the



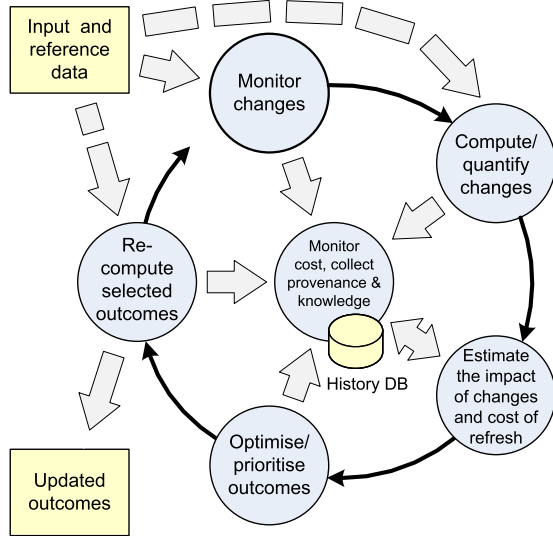
**Fig. 10.** Total time of the SVI workflow for a cohort of 33 patients depending on the approach taken to re-computation.

partial re-computation in all cases as there were not many changes in the columns relevant to SVI.

Fig. 10 shows the total re-computation time for the whole patient cohort including the time required to re-execute tasks with the difference sets and to run the partial re-computation with the complete new data when the impact function produced true. This figure emphasises the penalty for running the algorithm when the difference sets were large compared to actual new data. It also highlights the importance of the diff and impact functions. Clearly, the more accurate the functions are the higher runtime savings may be, which stems from two facts. Firstly, more accurate diff function tends to produce smaller difference sets which reduces time of task re-execution (cf. CV selective  $\delta$ -gen and  $\delta$ -SVI lines in Fig. 9). Secondly, more accurate impact function tends to produce false more frequently, and so the algorithm can more often avoid re-computation with the complete new version of the data (cf. the number of black squares vs the total number of patients affected by a change in Table 5).

**9. A blueprint for a generic and automated re-computation framework – challenges**

So far we have presented techniques that can be applied to reduce the cost of recurring re-computation, with reference to a



**Fig. 11.** The main loop in the ReComp framework that handles selective re-computation of the user process; thin black arrows denote the flow of control, thick arrows represent the flow of data.

single case study and without concern for the relative cost and benefits associated with the re-computation. Our long term goal is to generalise the approach into a reusable framework, which we call ReComp, that is able not only to carry out re-computations by automating a combination of the techniques we just illustrated, but also to help decision makers carry out a cost/benefit analysis to determine when selective re-computation is beneficial.

For this, ReComp must support a number of capabilities, above and beyond those just illustrated in Fig. 4. With reference to our execution model:

$$\langle y^t, c^t \rangle = \text{exec}(P, x^t, D^t) \quad (16)$$

these are:

1. Detect and quantify changes in input and reference data, i.e. by accepting data-specific  $\text{diff}_X()$  and  $\text{diff}_D()$  functions;
2. Estimate the impact of those changes on each member  $y^t \in Y$  in a population  $Y$  of prior outcomes, i.e. learn estimates of  $\text{diff}_Y(y^t, y^{t'})$  without having to compute  $y^{t'}$ , as well as estimates of the corresponding re-computation cost  $c^{t'}$ ;
3. Use the estimates to prioritise prior outcomes for re-computation, subject to a limited budget, and
4. Perform the re-computation of the corresponding instances of  $P$ , entirely or partially, as we have seen in this paper.

Note that, at this stage, we do not consider changes in  $P$  itself or any of its *software* dependencies (as opposed to the data dependencies). For simplicity we focus on changes in the data only and do not consider changes in the underlying processes. These are also relevant but require a separate formalisation, beyond the scope of this paper.

In practice, ReComp is configured as a *meta-process* that is able to (i) *monitor* instances of an underlying process  $P$  and record its provenance as well as details of its cost, (ii) *detect and quantify changes* in the data used by  $P$ , and (iii) *control* the re-execution of instances of  $P$ , on demand. These capabilities are summarised in the loop depicted in Fig. 11.

Central to ReComp is the idea that decisions about future re-executions are informed by analytics on the history of past executions. To make this possible, each execution of the form (16) (including re-executions) is controlled by ReComp, and generates metadata records that include:

- outcomes that are subject to revision;
- provenance of the outcome, either coarse-grained or fine-grained, depending on the underlying provenance recording facilities associated with the process runtime;
- execution cost, typically expressed as running time and data storage volume, again as detailed as allowed by the underlying system. For instance, our own WFMS, e-SC, provides block-level time recording and per-data-item storage, while other systems may only provide cumulative times.

Our long-term research hypothesis is that metadata analytics performed on such history database may yield viable models to estimate change impact and thus be able to prioritise re-computations vis-à-vis a limited budget. In the rest of this section we discuss a number of challenges that underpin the implementation of the ReComp framework.

### 9.1. Monitoring data changes

Managing multiple versions of large datasets is challenging. Firstly, observing changes in data usually requires source-specific solutions, as each resource is likely to expose a different version release mechanism – a version number being the simplest case. Secondly, the volume of data to be stored, multiplied by all the versions that might be needed for future re-computation, leads to prohibitively large storage requirements. Providers’ limitations in the versions they make available also translate into a challenge for ReComp, with some providers not offering access to different versions of their data at all.

A further issue is whether multiple changes to different data sources should be considered together or separately: in some cases it may be beneficial to group multiple changes to one resource instead of reacting immediately. For example, GeneMap updates are published daily, often with only a few rows changed. Thus, taking into account the cost of running the ReComp loop, it may be more effective to delay the loop and collect a number of updates, e.g. over a week.

### 9.2. Calculating and quantifying changes

Suppose two processes managed by ReComp retrieve different attributes from the same relational database  $D$ . Clearly, for each of these processes only changes to the relevant attributes matter. Thus, the  $\text{diff}()$  functions, such as those defined in Sec. 5, are not only type-specific but also query-specific. For  $n$  processes and  $m$  resources, this may potentially require  $n \cdot m$  specialised  $\text{diff}()$  functions. Whether we can find more effective ways to compute and measure data changes is an open question. Additionally, some input data may be unstructured or semi-structured and thus calculating the difference between two versions that is useful in the estimation of their impact may be challenging in itself.

### 9.3. Estimation impact and cost of refresh

We define the re-computation problem as finding the optimal selection of past invocation that can maximise the benefit of re-computation given changes in the input data and a budget constraint. Addressing this problem requires that we first learn impact estimators that can take into account the history of past executions, their cost and the changes in the input data, and can feed into the optimisation problem. This is a hard problem, however, which in particular involves estimating the difference between two outputs of process  $P$  given changes to some of its inputs. Clearly, some knowledge of the function that  $P$  implements is required, but that is also process-specific and so difficult to generalise into a reusable re-computation framework.



Recalling our example with SVI and ClinVar, we would like to predict whether or not a new variant added to the database will change patient diagnosis. The technique showed earlier allowed us to do so to some extent, as we were able to reduce the number of affected invocations from 495 to 71, yet more work is needed to find more accurate and more generic techniques.

The problem of learning cost estimators has been addressed in the recent past, but mainly for specific scenarios that are relevant to data analytics, namely workflow-based programming on clouds and grid [43,44]. But for instance [45] showed that runtime, especially in the case of machine learning algorithms, may depend on features that are specific to the input, and thus not easy to learn. That leaves the impact and cost estimation as an open challenge.

#### 9.4. Optimising the selection of past executions

Given a limited re-computation budget, and a measure of benefit of outcome refresh, we can address the further problem to select the past executions that are expected to maximise the benefit given the budget. Using the impact and cost estimators, we can formulate it as the 0–1 knapsack problem in which we want to find vector  $\mathbf{a} = [a_1 \dots a_n] \in \{0, 1\}^n$  that achieves:

$$\max \sum_{i=1}^n v_i a_i \quad \text{subject to} \quad \sum_{i=1}^n w_i a_i \leq C \quad (17)$$

where  $n$  is the number of past executions,  $v_i$  is the estimated change impact for execution  $i$ , and  $w_i$  is the estimated cost of its re-execution. Importantly, each data change event triggers an instance of (17) to be solved but due to expected high cost of re-computation it may be worth grouping a number of change events together. That adds complexity to the optimisation problem.

#### 9.5. Black box processes

Running the SVI example in the previous sections, we assumed that we have insight into the structure and semantics of process  $P$  managed by ReComp. That enabled us to effectively apply techniques for partial process re-execution. When  $P$  is a *black box* process, however, this is not possible and other techniques such as incremental computation [27,31,18] may be required. Regardless of the transparency of  $P$ , a common challenge is that for Big Data analytics intermediate data produced by the process (or memoised during incremental computation) often outgrow the actual inputs by orders of magnitude, and thus the cost of persisting all intermediate results may be prohibitive. An open problem, with some contribution from Woodman et al. [17], is to find techniques that could balance the choice of intermediate data to retain in view of a potential future re-computation, with its cost.

A separate challenge is that the actual re-execution of process  $P$  used in the past may not be straightforward. It may require re-deploying  $P$  on a new infrastructure and ensuring that the system and software dependencies are maintained correctly, or that the results obtained using new versions of third party libraries remain valid. Addressing these architectural and reproducibility issues is a research area of growing interest [23,21,22,46].

#### 9.6. History database

As mentioned, ReComp needs to collect and store both provenance and cost metadata. Recording cost requires the definition of a new format which, to the best of our knowledge, does not currently exist. Provenance, on the other hand, has been recorded using a number of formats, which are system-specific. Even when the PROV provenance model [19] is adopted, it can be used in different ways despite being designed to encourage interoperability.

Our recent study [47] shows that the ProvONE,<sup>12</sup> an extension to PROV, is a step forward to collect interoperable provenance traces, but is still limited as it assumes that the traced processes are similar and implemented as a workflow.

## 10. Conclusions and future work

Knowledge decay over time is an important issue that affects the value proposition of Big Data analytics. It is especially important for the next generation sequencing pipelines, in which algorithms and reference data continuously improve. As these pipelines require processing that can easily exceed hundreds of CPU-hours per patient cohort and as they become used on a wider scale,<sup>13</sup> relevant techniques to address knowledge decay and re-fresh pipeline results are required.

In this paper we presented our investigation into how selective re-computation can help address the knowledge decay issue. Using a case study in the area of clinical interpretation of genetic variants in humans, with a cohort of patients from the Institute of Genetic Medicine (IGM) at Newcastle University, we described three approaches to selective re-computation: at the process level (partial re-execution), data level (differential execution), and whole-cohort level (identification of scope of change).

Regarding partial re-execution, a special role is played by provenance, which we used to build the minimal process subgraph that requires re-execution. For differential execution, we used *diff()* functions to calculate difference sets between two versions of the input data and then, using these sets, to reduce the amount of processing needed. Finally, at the whole-cohort level we showed a significant reduction in the number of patient samples that required refresh. Overall, we were able to lower the cost of re-computation to about 10% of the total time needed for update the previous results. In the immediate future, our plan is to extend the study to a much larger cohort of over 1,500 patients [3], which will provide better figures on actual savings closer to real population scale.

Worth noting is that in this paper we discussed only the *lossless* approach to re-computation. Lossless re-computation is conservative in that any outcome on which the impact cannot be proved to be zero, regardless of how small, will be refreshed. In contrast, *lossy re-computation* also seeks to reduce the amount of work performed on previously computed outcome. However, lossy re-computation would try to quantitatively *estimate* the extent of the impact, and use the estimates to decide whether and when to refresh the outcomes. We view this as a more general re-computation decision problem which involves a cost/benefit analysis.

This study on variant interpretation informs the more ambitious ReComp project.<sup>14</sup> Our immediate next step is to apply the techniques presented in this paper to other parts of the variant calling pipeline. We are now developing a generic meta-process that can observe changes and control re-computation for a variety of underlying, resource-intensive analytics processes, as well as support business-level re-computation decisions vis-à-vis a resource budget. In Sec. 9 we outlined a number of the research and technical challenges associated with this vision.

## Acknowledgements

Funding: This work has been supported by EPSRC in the UK [grant number EP/N01426X/1]; and a grant from the Microsoft Azure for Research programme.

<sup>12</sup> <https://purl.dataone.org/provone-v1-dev>.

<sup>13</sup> <https://www.genomicsengland.co.uk/the-100000-genomes-project/>.

<sup>14</sup> <http://recomp.org.uk>.

Appendix A. Input data

Table A.6 Basic properties of a set of patient variant files used in the experiments.

Table with 4 columns: Phenotype hypothesis, Variant file, Record count, File size [MB]. Rows include Alzheimer's disease, CADASIL, and Frontotemporal dementia - Amyotrophic lateral sclerosis with various variant files and their respective counts and sizes.

Table A.7 Basic properties of the OMIM GeneMap and ClinVar reference databases used in the experiments.

Table with 4 columns: Database, Version, Record count, File size [MB]. Rows include OMIM GeneMap and NCBI ClinVar with their respective versions, record counts, and file sizes.

Appendix B. Supplementary material

Supplementary material related to this article can be found online at https://doi.org/10.1016/j.bdr.2018.06.001.

Table A.8 Changes observed in the output of the SVI tool when executed with the difference sets computed for NCBI ClinVar reference database using the generic delta function; ■ denotes the need for re-execution with the complete new version of ClinVar (D\_ac != empty or D\_r != empty), \* denotes only task re-execution with the difference sets (D\_ac = empty and D\_r = empty).

Table with 16 columns for variant files (D\_1071 to B\_0198) and 11 rows for ClinVar versions (08/15 to 10/16). The table uses a grid of symbols (■, \*, .) to indicate changes in the output of the SVI tool for each combination of variant file and ClinVar version.

## References

- [1] R. Do, S. Kathiresan, G.R. Abecasis, Exome sequencing and complex disease: practical aspects of rare variant association studies, *Hum. Mol. Genet.* 21 (R1) (2012) R1–R9, <https://doi.org/10.1093/hmg/dds387>.
- [2] H. Buermans, J. den Dunnen, Next generation sequencing technology: advances and applications, *Biochim. Biophys. Acta, Mol. Basis Dis.* 1842 (10) (2014) 1932–1941, <https://doi.org/10.1016/j.bbadis.2014.06.015>.
- [3] M.J. Keogh, W. Wei, I. Wilson, J. Coxhead, S. Ryan, S. Rollinson, H. Griffin, M. Kurzawa-Akanbi, M. Santibanez-Koref, K. Talbot, M.R. Turner, C.-A. McKenzie, C. Troakes, J. Attems, C. Smith, S. Al Sarraj, C.M. Morris, O. Ansoorge, S. Pickering-Brown, J.W. Ironside, P.F. Chinnery, Genetic compendium of 1511 human brains available through the UK Medical Research Council Brain Banks Network Resource, *Genome Res.* 27 (1) (2017) 165–173, <https://doi.org/10.1101/gr.210609.116>.
- [4] J. Cała, E. Marei, Y. Xu, K. Takeda, P. Missier, Scalable and efficient whole-exome data processing using workflows on the cloud, *Future Gener. Comput. Syst.* 65 (2016) 153–168, <https://doi.org/10.1016/j.future.2016.01.001>, special Issue on Big Data in the Cloud.
- [5] Y. Qin, H.K. Yalamanchili, J. Qin, B. Yan, J. Wang, The current status and challenges in computational analysis of genomic big data, *Big Data Res.* 2 (1) (2015) 12–18, <https://doi.org/10.1016/j.bdr.2015.02.005>.
- [6] L. Hood, S.H. Friend, Predictive, personalized, preventive, participatory (P4) cancer medicine, *Nat. Rev. Clin. Oncol.* 8 (3) (2011) 184, <https://doi.org/10.1038/nrclinonc.2010.227>.
- [7] P. Missier, E. Wijaya, R. Kirby, M. Keogh, SVI: a simple single-nucleotide human variant interpretation tool for clinical use, in: N. Ashish, J.-L. Ambite (Eds.), *Data Integration in the Life Sciences*, Springer International Publishing, 2015, pp. 180–194.
- [8] E. Garrison, G. Marth, Haplotype-based variant detection from short-read sequencing, [arXiv:1207.3907](https://arxiv.org/abs/1207.3907).
- [9] S. Sandmann, A.O. de Graaf, M. Karimi, B.A. van der Reijden, E. Hellström-Lindberg, J.H. Jansen, M. Dugas, Evaluating variant calling tools for non-matched next-generation sequencing data, *Sci. Rep.* 7 (2017) 43169, <https://doi.org/10.1038/srep43169>, EP.
- [10] D.M. Church, V.A. Schneider, K.M. Steinberg, M.C. Schatz, A.R. Quinlan, C.-S. Chin, P.A. Kitts, B. Aken, G.T. Marth, M.M. Hoffman, J. Herrero, M.L.Z. Mendoza, R. Durbin, P. Flicek, Extending reference assembly models, *Genome Biol.* 16 (1) (2015) 13, <https://doi.org/10.1186/s13059-015-0587-3>.
- [11] V.A. Schneider, T. Graves-Lindsay, K. Howe, N. Bouk, H.-C. Chen, P.A. Kitts, T.D. Murphy, K.D. Pruitt, F. Thibaud-Nissen, D. Albracht, R.S. Fulton, M. Kremitzki, V. Magrini, C. Markovic, S. McGrath, K.M. Steinberg, K. Auger, W. Chow, J. Collins, G. Harden, T. Hubbard, S. Pelan, J.T. Simpson, G. Threadgold, J. Torrance, J.M. Wood, L. Clarke, S. Koren, M. Boitano, P. Peluso, H. Li, C.-S. Chin, A.M. Phillippy, R. Durbin, R.K. Wilson, P. Flicek, E.E. Eichler, D.M. Church, Evaluation of grch38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly, *Genome Res.* 27 (5) (2017) 849–864, <https://doi.org/10.1101/gr.213611.116>, <http://genome.cshlp.org/content/27/5/849.full.pdf+html>, <http://genome.cshlp.org/content/27/5/849.abstract>.
- [12] S. Hwang, E. Kim, I. Lee, E.M. Marcotte, Systematic comparison of variant calling pipelines using gold standard personal exome variants, *Sci. Rep.* 5 (December) (2015) 17875, <https://doi.org/10.1038/srep17875>.
- [13] E.T. Cirulli, D.B. Goldstein, Uncovering the roles of rare variants in common disease through whole-genome sequencing, *Nature reviews, Genetics* 11 (6) (2010) 415–425, <https://doi.org/10.1038/nrg2779>.
- [14] I. Altintas, O. Barney, E. Jaeger-Frank, Provenance collection support in the Kepler Scientific Workflow System, *Work* 4145 (2006) 118–132, [https://doi.org/10.1007/11890850\\_14](https://doi.org/10.1007/11890850_14).
- [15] H. Lakhani, R. Tahir, A. Aqil, F. Zaffar, D. Tariq, A. Gehani, Optimized rollback and re-computation, in: 2013 46th Hawaii International Conference on System Sciences, IEEE, 2013, pp. 4930–4937.
- [16] Yaxiong Zhao, Jie Wu, Cong Liu, Dache: a data aware caching for big-data applications using the MapReduce framework, *Tsinghua Sci. Technol.* 19 (1) (2014) 39–50, <https://doi.org/10.1109/TST.2014.6733207>.
- [17] S. Woodman, H. Hiden, P. Watson, Workflow provenance: an analysis of long term storage costs, in: Proceedings of the 10th Workshop on Workflows in Support of Large-Scale Science, 2015, pp. 9:1–9:9.
- [18] F.D. McSherry, D.G. Murray, R. Isaacs, M. Isard, Differential dataflow, in: 6th Biennial Conference on Innovative Data Systems Research, CIDR '13, 2013, [http://cidrdb.org/cidr2013/Papers/CIDR13\\_Paper111.pdf](http://cidrdb.org/cidr2013/Papers/CIDR13_Paper111.pdf).
- [19] PROV-DM: The PROV Data Model, Technical Report, World Wide Web Consortium, Apr. 2013, <https://www.w3.org/TR/prov-dm/>.
- [20] J.F. Pimentel, J. Freire, V. Braganholo, L. Murta, Tracking and analyzing the evolution of provenance from scripts, in: Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 9672, Springer International Publishing, 2016, pp. 16–28.
- [21] J. Freire, N. Fuhr, A. Rauber, Reproducibility of data-oriented experiments in e-science, in: Dagstuhl Seminar 16041, vol. 6, 2016.
- [22] L.C. Burgess, D. Crotty, D. de Roure, J. Gibbons, C. Goble, P. Missier, R. Mortier, T.E. Nichols, R. O'Beirne, Alan Turing Institute Symposium on Reproducibility for Data-Intensive Research – Final Report, Tech. Rep., 2016.
- [23] V. Stodden, F. Leisch, R.D. Peng, Implementing Reproducible Research, Chapman & Hall/CRC The R Series, CRC Press, 2014.
- [24] M. Herschel, R. Diestelkämper, H. Ben Lahmar, A survey on provenance: what for? What form? What from?, *Vldb J.* 26 (6) (2017) 1–26, <https://doi.org/10.1007/s00778-017-0486-1>.
- [25] B. Ludäscher, I. Altintas, C. Berkley, D. Higgins, E. Jaeger, M. Jones, E.A. Lee, J. Tao, Y. Zhao, Scientific workflow management and the Kepler system, *Concurr. Comput., Pract. Exp.* 18 (10) (2005) 1039–1065, <https://doi.org/10.1002/cpe.994>.
- [26] L. Bavoil, S. Callahan, P. Crossno, J. Freire, C. Scheidegger, C. Silva, H. Vo, Vis-Trails: enabling interactive multiple-view visualizations, in: IEEE Visualization, 2005., no. Dx, VIS 05, IEEE, 2005, pp. 135–142.
- [27] U.A. Acar, G.E. Blueloch, M. Blume, R. Harper, K. Tangwongsan, An experimental analysis of self-adjusting computation, *ACM Trans. Program. Lang. Syst.* 32 (1) (2009) 1–53, <https://doi.org/10.1145/1596527.1596530>.
- [28] G. Ramalingam, T. Reps, A categorized bibliography on incremental computation, in: Proceedings of the 20th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '93, 1993, pp. 502–510.
- [29] L. Popa, M. Budiu, Y. Yu, M. Isard, DryadInc: reusing work in large-scale computations, in: HotCloud'09 Workshop on Hot Topics in Cloud Computing, 2009, pp. 2–6, [http://static.usenix.org/events/hotcloud09/tech/full\\_papers/popa.pdf](http://static.usenix.org/events/hotcloud09/tech/full_papers/popa.pdf).
- [30] Y. Bu, B. Howe, M. Balazinska, M.D. Ernst HaLoop, *Proc. VLDB Endow.* 3 (1–2) (2010) 285–296, <https://doi.org/10.14778/1920841.1920881>.
- [31] P. Bhatotia, A. Wieder, R. Rodrigues, U.A. Acar, R. Pasquin, Incoop: MapReduce for incremental computations, in: Proceedings of the 2nd ACM Symposium on Cloud Computing, SOCC '11, 2011, pp. 1–14.
- [32] A.G. Bin Saadon, H.M.O. Mokhtar, liHadoop: an asynchronous distributed framework for incremental iterative computations, *J. Big Data* 4 (1) (2017) 24, <https://doi.org/10.1186/s40537-017-0086-3>.
- [33] P. Bhatotia, P. Fonseca, U.A. Acar, B.B. Brandenburg, R. Rodrigues iThreads, *ACM SIGARCH Comput. Archit. News* 43 (1) (2015) 645–659, <https://doi.org/10.1145/2786763.2694371>.
- [34] H. Hiden, S. Woodman, P. Watson, J. Cała, Developing cloud applications using the e-science central platform, *Philos. Trans. R. Soc., Math. Phys. Eng. Sci.* 371 (1983), <https://doi.org/10.1098/rsta.2012.0085>.
- [35] V. Cuevas-Vicenttin, B. Ludäscher, P. Missier, K. Belhajjame, F. Chirigati, Y. Wei, S. Dey, P. Kianmajid, D. Koop, S. Bowers, I. Altintas, C. Jones, M.B. Jones, L. Walker, P. Slaughter, B. Leinfelder, Y. Cao, ProvONE: a PROV Extension Data Model for Scientific Workflow Provenance, Technical Report, DataONE Cyberinfrastructure Working Group, May 2016, <http://jenkins-1.dataone.org/jenkins/view/DocumentationProjects/job/ProvONE-Documentation-trunk/ws/provenance/ProvONE/v1/provone.html>.
- [36] L. Moreau, P. Missier, J. Cheney, S. Soiland-Reyes, PROV-n: The Provenance Notation, Tech. Rep., 2012, <http://www.w3.org/TR/prov-n/>.
- [37] Y. Chen, U.A. Acar, K. Tangwongsan, Functional programming for dynamic and large data with self-adjusting computation, in: Proceedings of the 19th ACM SIGPLAN International Conference on Functional Programming, ICFP '14, 2014, pp. 227–240.
- [38] Y. Cui, J. Widom, Lineage tracing for general data warehouse transformations, *Vldb J.* 12 (1) (2003) 41–58, <https://doi.org/10.1007/s00778-002-0083-8>.
- [39] P. Missier, N.W. Paton, K. Belhajjame, Fine-grained and efficient lineage querying of collection-based workflow provenance, in: Proceedings of the 13th International Conference on Extending Database Technology, EDBT '10, 2010, p. 299.
- [40] Z. Zhang, E.R. Sparks, M.J. Franklin, Diagnosing machine learning pipelines with fine-grained lineage, in: Proceedings of the 26th International Symposium on High-Performance Parallel and Distributed Computing, HPDC '17, ACM Press, New York, USA, 2017, pp. 143–153.
- [41] Y. Cui, J. Widom, Practical lineage tracing in data warehouses, in: Proceedings of the 16th International Conference on Data Engineering, 2000, pp. 367–378.
- [42] P.-A. Larson, J. Zhou, Efficient maintenance of materialized outer-join views, in: 2007 IEEE 23rd International Conference on Data Engineering, IEEE, 2007, pp. 56–65.
- [43] I. Pietri, G. Juve, E. Deelman, R. Sakellariou, A performance model to estimate execution time of scientific workflows on the cloud, in: 2014 9th Workshop on Workflows in Support of Large-Scale Science, IEEE, 2014, pp. 11–19.
- [44] M.J. Malik, T. Fahringer, R. Prodan, Execution time prediction for grid infrastructures based on runtime provenance data, in: Proceedings of the 8th Workshop on Workflows in Support of Large-Scale Science, WORKS '13, ACM Press, New York, New York, USA, 2013, pp. 48–57.
- [45] T. Miu, P. Missier, Predicting the execution time of workflow activities based on their input features, in: Proceedings – 2012 SC Companion: High Performance Computing, Networking Storage and Analysis, SCC 2012, 2012, pp. 64–72.
- [46] R. Qasha, J. Cała, P. Watson, A framework for scientific workflow reproducibility in the cloud, in: 2016 IEEE 12th International Conference on e-Science, e-Science, IEEE, 2016, pp. 81–90.
- [47] W. Oliveira, P. Missier, K. Ocaña, D. de Oliveira, V. Braganholo, Analyzing provenance across heterogeneous provenance graphs, in: Ipaw, vol. 5272, 2016, pp. 57–70.