**Use Case Oncology Data Management**

Dr. Matthieu-P. Schapranow

Data Management for Digital Health
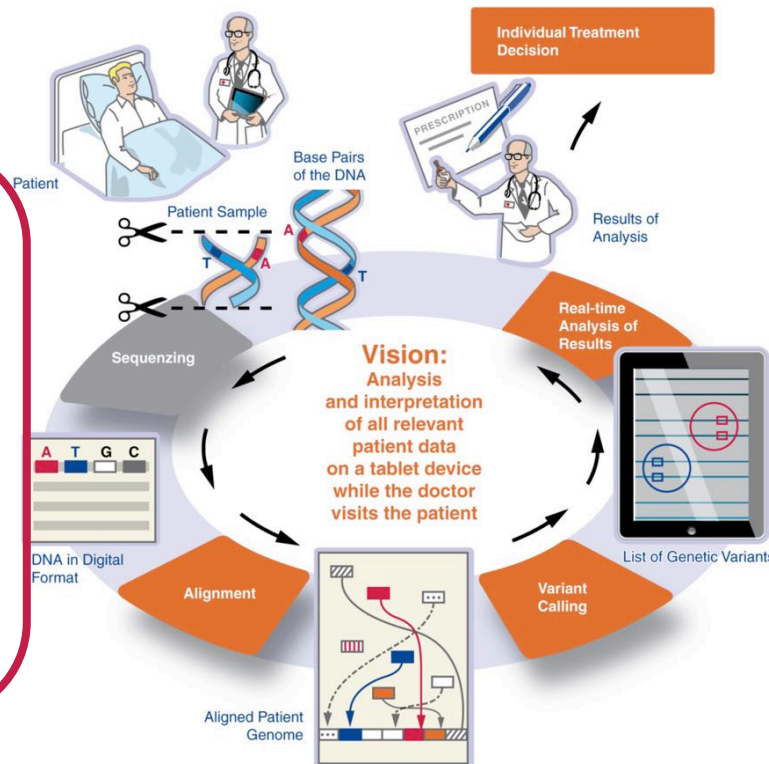
Summer 2017

# From Raw Genome Data to Analysis

- **Sequencing**: Acquire digital DNA data

- **Alignment**: Reconstruction of complete genome with snippets

- **Variant Calling**: Identification of genetic variants

- **Data Annotation**: Linking genetic variants with research findings

# Genome Data Processing Pipelines
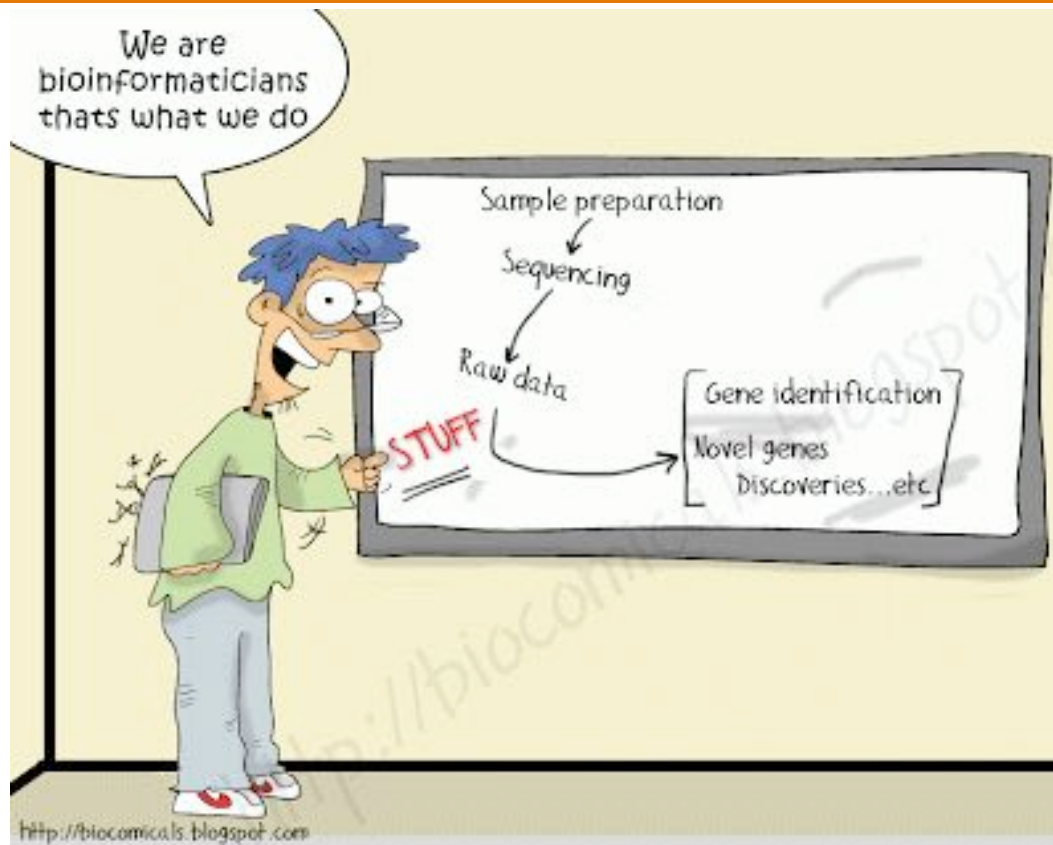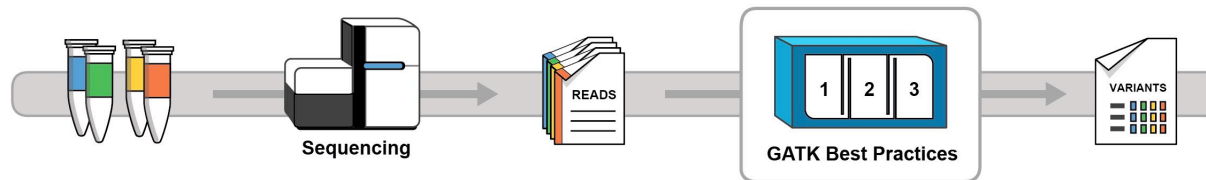## Motivation

- Currently Genome Data Processing Pipelines (GDPPs) consist of individually combined command line scripts

- Complex to
  - Understand,
  - Maintain and update, and
  - Reproduce

- Objective: Model GDPPs in a...
  - Graphical and machine-readable representation
  - Reproducible and exchangeable format

Use Case Oncology
Data Management

Data Management for
Digital Health, Summer
2017
4

- The BROAD is a joint institute of MIT and Harvard established 2004 in Cambridge, MA

- Genome Analysis Toolkit (GATK) focuses on variant detection

- Open-source tools and shared best-practices

☑ **GATK Best Practices**

Lots of workflows that people call Best Practices really aren't.

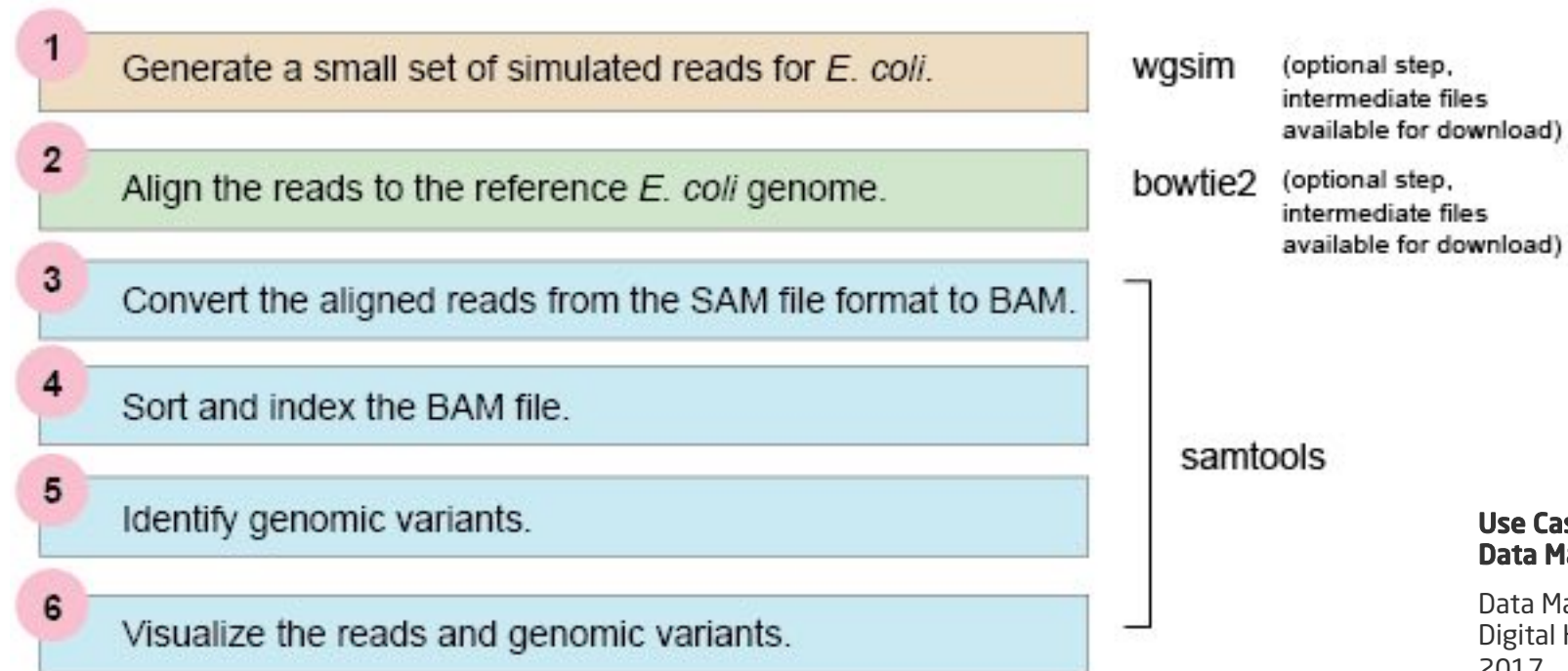https://software.broadinstitute.org/gatk/best-practices/

# Genome Data Processing Pipelines
## State of the Art

*bwa aln hg19.fa sample.fastq | bwa samse hg19.fa – sample.fastq | samtools view -Su - | samtools sort …*

- Concatenation of command line tools reading/writing files

- Lack of standardization and exchangeability

- Requires dedicated expertise for

  □ Setup and configuration,

  □ Error handling, and

  □ Scalable processing

- Objective: Enable modeling, configuration, and execution of reproducible pipelines without involving IT experts

**Use Case Oncology Data Management**

Data Management for Digital Health, Summer 2017

6

# Let us Make it a Workflow!



| | | |
|---|---|---|
| 1 | Generate a small set of simulated reads for *E. coli*. | wgsim (optional step, intermediate files available for download) |
| 2 | Align the reads to the reference *E. coli* genome. | bowtie2 (optional step, intermediate files available for download) |
| 3 | Convert the aligned reads from the SAM file format to BAM. | samtools |
| 4 | Sort and index the BAM file. | |
| 5 | Identify genomic variants. | |
| 6 | Visualize the reads and genomic variants. | |

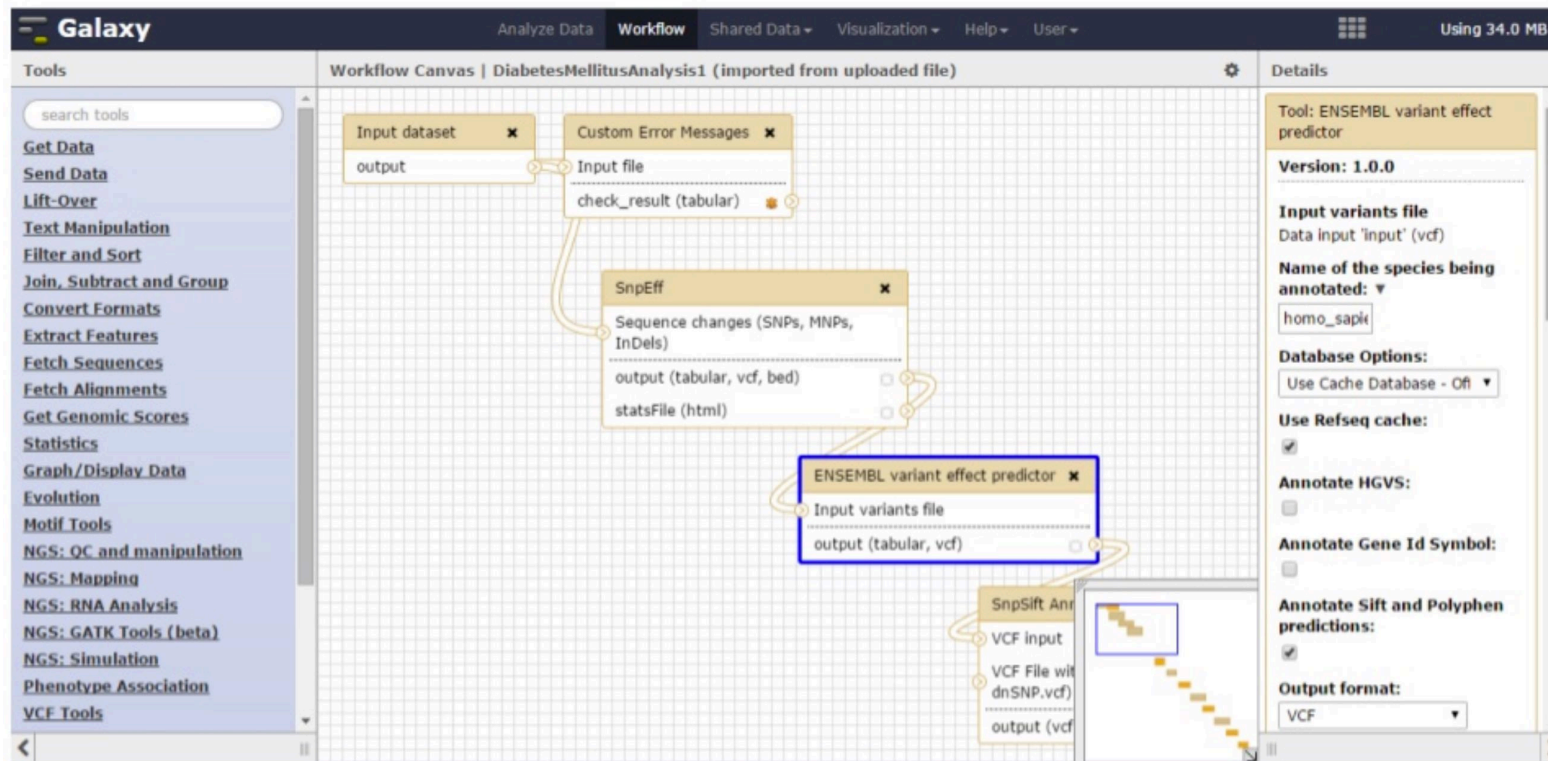http://biobits.org/samtools_primer.html

# Galaxy
# Workbench

- Open-source, web-based platform

- Supports data-intensive research

- Focuses on process automation and high-throughput sequencing

# Galaxy
# Workflow Modeling

https://usegalaxy.org/

# DKFZ One Touch Platform

- IT process automation at DKFZ, HD
- Builds upon OpenStack and Vagrant to reduce setup time
- Workflow managed by SeqWare Pipeline Manager
- Special-purpose developed for DKFZ requirements



https://seqware.github.io/docs/6-pipeline/

# Google Genomics

- Integration of existing Google services to genome data processing

**Use Case Oncology Data Management**

Data Management for Digital Health, Summer 2017

11

# Business Process Modeling and Notation (BPMN) 2.0

- Used for functional modeling of business processes and workflows

- Graphical notation addresses business and technical users → intuitive modeling and understanding

- Can be serialized and exchanged using XML Process Definition Language (XPDL)

```
<?xml version="1.0" encoding="UTF-8"?>
<zdef-2030967014:Package xmlns="" xmlns:xpdExt="http://www.tibco.com/XPD/xpdExtens
  <zdef-2030967014:ConformanceClass GraphConformance="NON-BLOCKED" BPMNModelPortab
  <zdef-2030967014:Script Type="http://www.w3.org/1999/XPath"/>
  <Pools xmlns="http://www.wfmc.org/2008/XPDL2.1">
    <Pool BoundaryVisible="false" MainPool="true" Process="MainPool-process" Orier
      <NodeGraphicsInfos>
        <NodeGraphicsInfo FillColor="#ffffff" Height="0.0" Width="0.0" BorderColc
          <Coordinates XCoordinate="0.0" YCoordinate="0.0"/>
      </NodeGraphicsInfo>
      </NodeGraphicsInfos>
    </Pool>
  </Pools>
  <WorkflowProcesses xmlns="http://www.wfmc.org/2008/XPDL2.1">
    <WorkflowProcess AdhocOrdering="Sequential" ProcessType="None" Status="None" S
      <ActivitySets>
        <ActivitySet AdHocOrdering="Sequential" Id="sid-A846876F-9749-41F9-93DE-6C
        <ActivitySet AdHocOrdering="Sequential" Id="sid-10D16CD8-AAEF-4694-A5C4-75
      </ActivitySets>
```
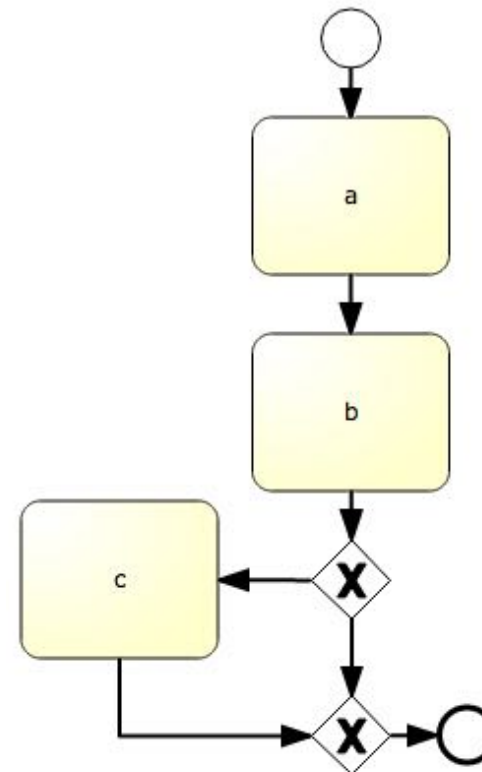
**Use Case Oncology Data Management**

Data Management for Digital Health, Summer 2017

12

# BPMN 2.0: Basic Notation Overview

# Reproducibility
## Modeling of Data Analysis Pipelines

1. Design time (researcher, process expert)
   - □ Definition of parameterized process model
   - □ Uses graphical editor and jobs from repository

2. Configuration time (researcher, lab assistant)
   - □ Select model and specify parameters, e.g. aln opts
   - □ Results in model instance stored in repository

3. Execution time (researcher)
   - □ Select model instance
   - □ Specify execution parameters, e.g. input files

# Standardized Graphical Modeling

- Graphical modeling notation compliant with Business Process Modeling and Notation 2.0 extended by

  □ Modular structure

  □ Parallelization annotations

  □ Parameters and variables

- Model descriptions are stored within IMDB

- Model instances are transformed into graph structure for execution by dedicated runtime environment

# BPMN Example

Data Management for
Digital Health, Summer
2017

16

# Persisting Pipelines
# XML Process Definition Language

```xml
<xpdl:Activity CompletionQuantity="1" Id="newpkg1_wp1_act2" Name="BWA"
    <xpdl:Implementation>
        <xpdl:No/>
    </xpdl:Implementation>
    <xpdl:Performers>
        <xpdl:Performer>newpkg1_wp1_par1</xpdl:Performer>
    </xpdl:Performers>
    <xpdl:NodeGraphicsInfos>
        <xpdl:NodeGraphicsInfo BorderColor="#000000" FillColor="#99FF99
            <xpdl:Coordinates XCoordinate="239.0" YCoordinate="219.0"/>
        </xpdl:NodeGraphicsInfo>
    </xpdl:NodeGraphicsInfos>
</xpdl:Activity>
```
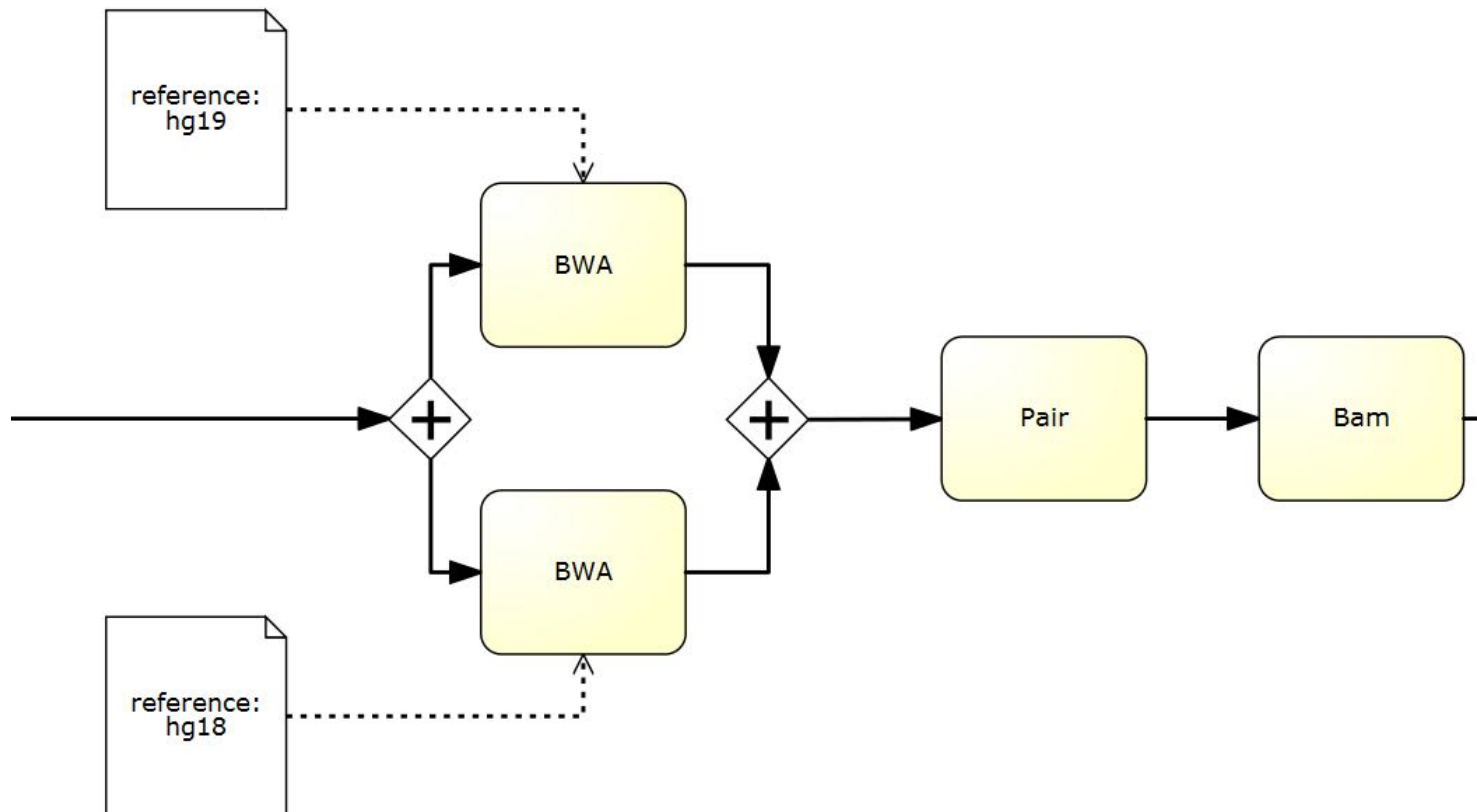
```xml
<xpdl:Artifacts>
    <xpdl:Artifact ArtifactType="DataObject" Id="newpkg1_1" Name="newpkg1_1">
        <xpdl:DataObject Id="newpkg1_1" Name="reference:hg19"/>
        <xpdl:NodeGraphicsInfos>
            <xpdl:NodeGraphicsInfo BorderColor="#000000" FillColor="#E8EEF7"
                <xpdl:Coordinates XCoordinate="239.0" YCoordinate="74.0"/>
            </xpdl:NodeGraphicsInfo>
        </xpdl:NodeGraphicsInfos>
    </xpdl:Artifact>
```
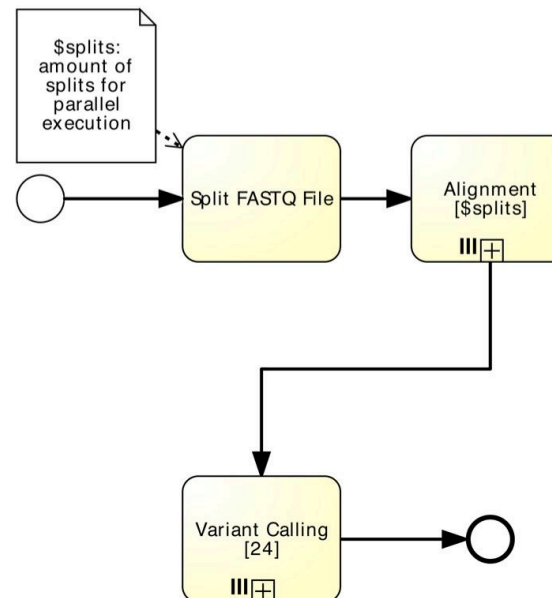
# Database Structure

## PIPELINES.MODELS

| ID | NAME | CONTENT | TYPE |
|----|------|---------|------|
| 154 | BWA2 | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 150 | VCFImport | <?xml version="1.0" encoding="UTF-8"?> | |
| 146 | Alignme... | <?xml version="1.0" encoding="UTF-8"?> | |
| 144 | Bowtie2_s | <?xml version="1.0" encoding="UTF-8"?> | alignme... |
| 143 | TMAP4 | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 141 | TMAP2 | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 140 | TMAP1 | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 139 | HANA_Al... | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 138 | BWA | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 137 | Bowtie2 | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 136 | Bowtie | <?xml version="1.0" encoding="UTF-8"?> | alignment |
| 135 | HANA_Al... | <?xml version="1.0" encoding="UTF-8"?> | alignme... |
| 134 | BWA_s | <?xml version="1.0" encoding="UTF-8"?> | alignme... |
| 133 | Bowtie_s | <?xml version="1.0" encoding="UTF-8"?> | alignme... |
| 129 | Optimized | <?xml version="1.0" encoding="UTF-8"?> | main |
| 128 | Standard | <?xml version="1.0" encoding="UTF-8"?> | main |
| 102 | Paired_A... | <?xml version="1.0" encoding="UTF-8"?> | main |

## PIPELINES.PIPELINES

| ID | NAME | MODE | CONFIG | SUBTASKS |
|----|------|------|--------|----------|
| 106 | Test | 129 | {"split_count":"25","reference":"hg19"} | {"alignment_speed":"135"} |
| 105 | Test5 | 156 | {"reference":"hg19"} | {} |
| 104 | Stand... | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"154"} |
| 103 | test | 153 | {"reference":""} | {} |
| 78 | TMAP4 | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"143"} |
| 77 | TMAP3 | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"142"} |
| 76 | TMAP2 | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"141"} |
| 74 | HANA... | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"139"} |
| 73 | BWA | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"138"} |
| 72 | Bowtie2 | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"137"} |
| 71 | Bowtie | 128 | {"split_count":"3","reference":"hg19"} | {"alignment":"136"} |
| 70 | Bowtie... | 129 | {"split_count":"3","reference":"hg19"} | {"alignment_speed":"144"} |
| 69 | HANA... | 129 | {"split_count":"3","reference":"hg19"} | {"alignment_speed":"135"} |
| 68 | BWA_s | 129 | {"split_count":"25","reference":"hg19"} | {"alignment_speed":"134"} |
| 67 | Bowtie_s | 129 | {"split_count":"3","reference":"hg19"} | {"alignment_speed":"133"} |

# Traditional vs. IMDB Approach

- Processing and results are kept within IMDB
- Optimization reduced execution time by >**50%**
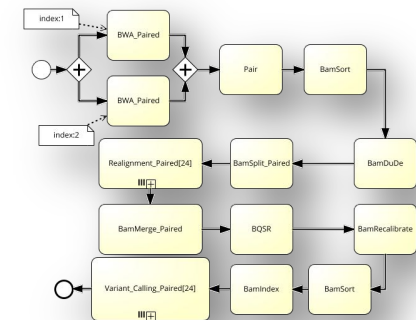
# Available Tools

- Pipeline categories (each traditional and optimized):

  - Single read,

  - Paired read,

  - Amplicon

- Alignment: IMDB-based, TMAP, BWA, Bowtie, Isaac

- Variant Calling: IMDB-based, BcfTools, GATK

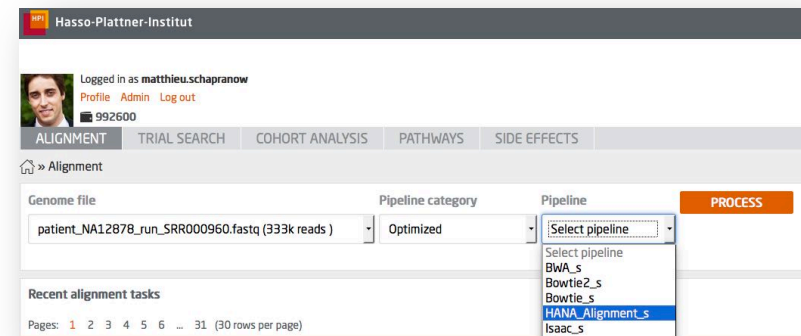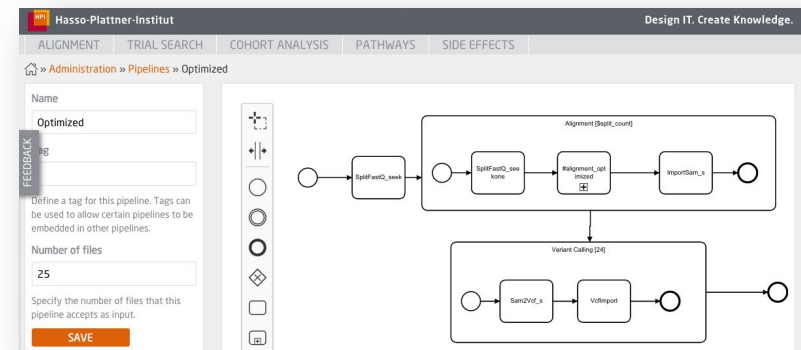- Intermediate steps: SAMtools, Picard, GATK

# From Model to Execution

1. **Design time** (researcher, process expert)
   - ☐ Definition of parameterized process model
   - ☐ Uses graphical editor and jobs from repository

2. **Configuration time** (researcher, lab assistant)
   - ☐ Select model and specify parameters, e.g. aln opts
   - ☐ Results in model instance stored in repository

3. **Execution time** (researcher)
   - ☐ Select model instance
   - ☐ Specify execution parameters, e.g. input files

# Execution of GDPPs

- Uses workflow, which is...
  - Predefined by a subject-matter expert
  - Preconfigured for a specific run or set of experiments
- Requires only minimal configuration whilst enabling reproducibility

**New alignment task**

(1) Choose pipeline | **Configure execution** | Select file(s)

Pipeline #alignment_speed
BWA_s

Variable $split_count
10

Variable $reference
hg19

**SELECT FILES >**

**New alignment task**

(2) Choose pipeline | Configure execution | **Select file(s)**

Genome file #1
Choose a file | or upload a new one

Choose a file
*User files*
CMV2_d15_cDNA_hTCRAlpha_454.fastq (6k reads)
CMV2_d9_cDNA_hTCRAlpha_454.fastq (44k reads)
CMV2_enriched_cDNA_hTCRBeta_454.fastq (8k reads)
ERR005584.filt.fastq (65k reads)
ERR031969.filt.fastq (107k reads)
ERR047877.filt.fastq (46k reads)
HN-10927_S13_L001_001_1.fastq (582k reads)
HN-10980A_S14_L001_001_1.fastq (276k reads)
HN-10980A_S14_L001_001_2.fastq (276k reads)
HN-10980A_S14_L001_R1_001_1.fastq (276k reads)
Sezary2_PB_cDNA_hTCRBeta_454.fastq (77k reads)
Sezary7_PB_cDNA_hTCRAlpha_454.fastq (47k reads)
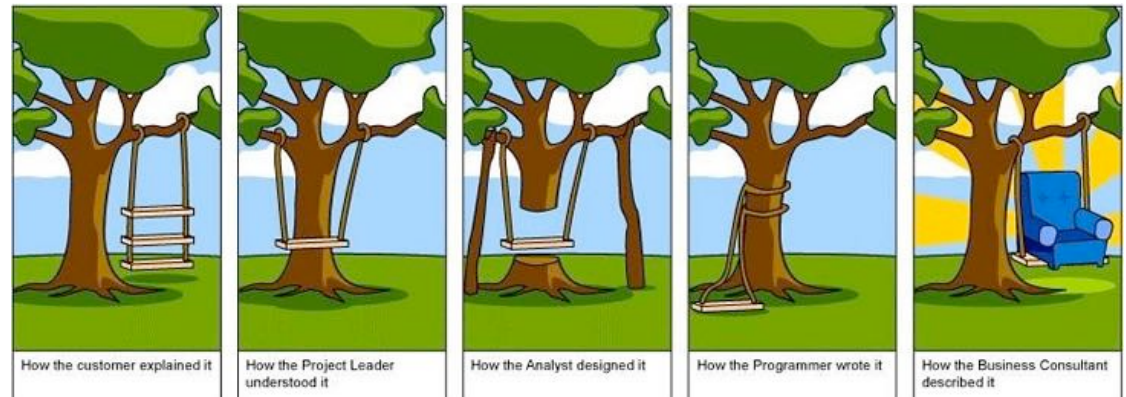
# Execution Environment for GDPPs
# Software Requirements in Life Sciences

- Requirements
  - Managed services
  - Reproducibility
  - Real-time data analysis of big data
- Restrictions
  - Data privacy
  - Data locality
  - Volume of big medical data

- Solution?
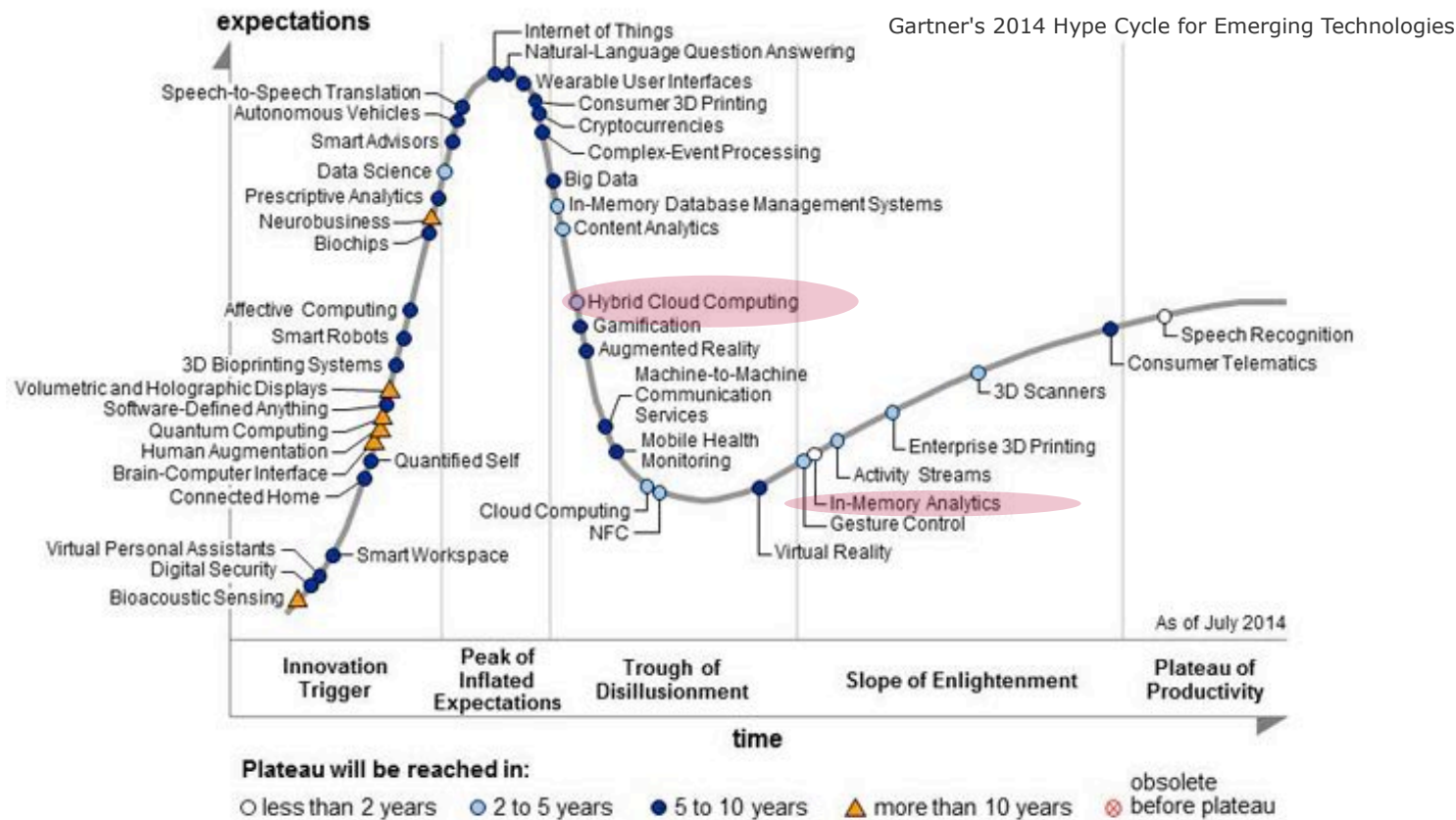  - Federated In-Memory Database System vs. Cloud Computing



http://stevedempsen.blogspot.de/2013/08/agile-software-requirements-comic.html

**Use Case Oncology
Data Management**

Data Management for
Digital Health, Summer
2017
23

# Execution Environment for GDPPs
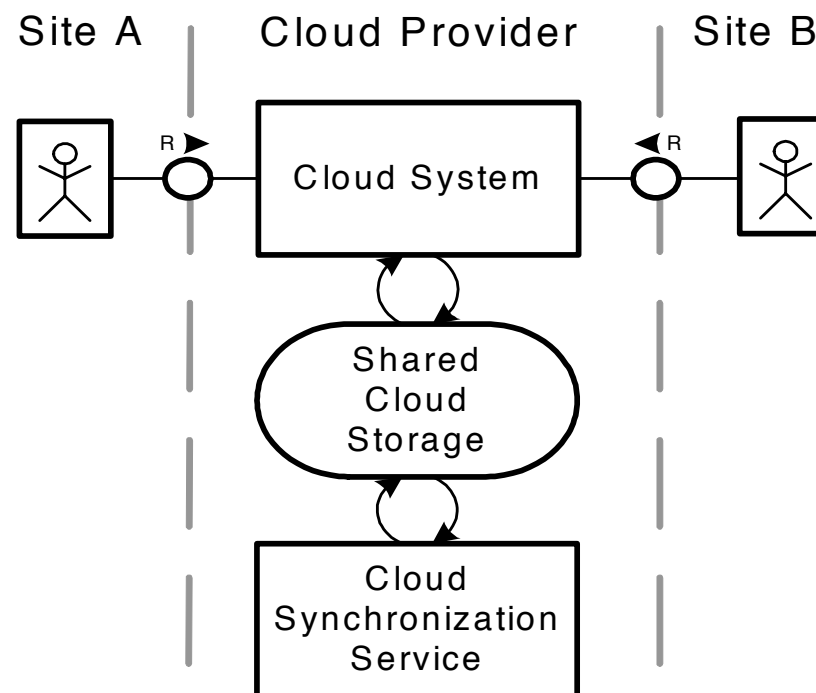# Cloud vs. On-Premise?



Gartner's 2014 Hype Cycle for Emerging Technologies

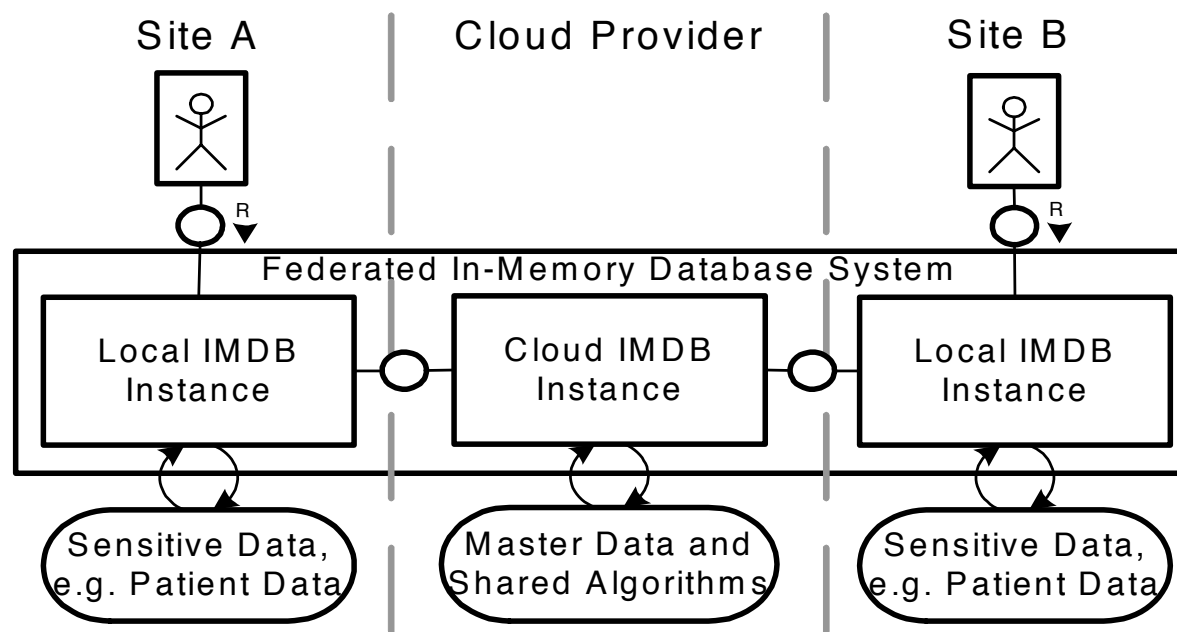# Multiple Cloud Service Providers

# A Single Service Provider

# Multiple Sites Forming the
# Federated In-Memory Database System



Site A     Cloud Provider     Site B

Federated In-Memory Database System

Local IMDB Instance     Cloud IMDB Instance     Local IMDB Instance

Sensitive Data, e.g. Patient Data     Master Data and Shared Algorithms     Sensitive Data, e.g. Patient Data

# Network Setup
# Site-to-Site VPN

Hasso Plattner Institut

LAN Site A
141.80.177.0/23

Site-to-Site VPN Tunnel

LAN Site B
192.168.10.0/24

**Public Internet**

VPN Gateway

VPN Gateway

**MDC**

**Consumer**

**HPI**

**Managed Services Provider**

**Use Case Oncology Data Management**

Data Management for Digital Health, Summer 2017

28

# Federated In-Memory Database (FIMDB)
## Incorporating Local Compute Resources



Federated In-Memory
Database Instance,
Algorithms, and
Applications Managed
by Service Provider

Federated In-Memory
Database Instances

Master Data
Managed by
Service Provider

Sensitive Data
reside at Site

**Use Case Oncology
Data Management**

Data Management for
Digital Health, Summer
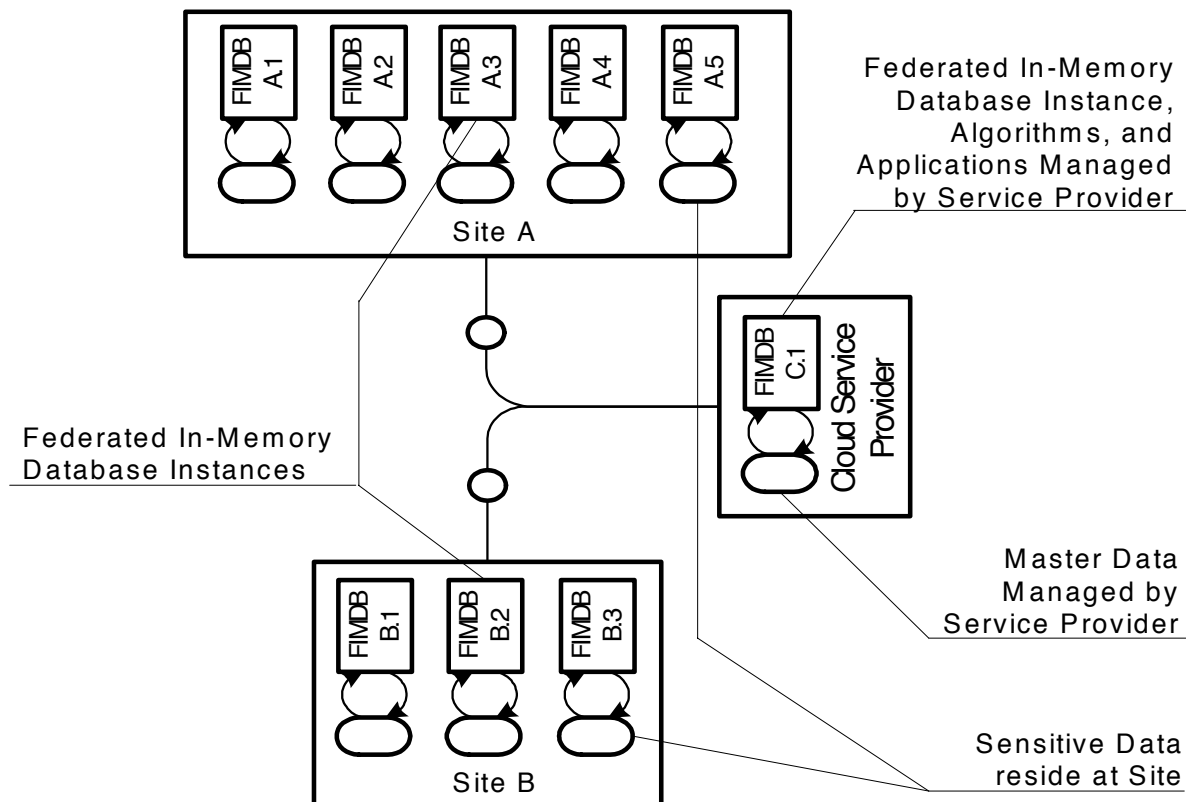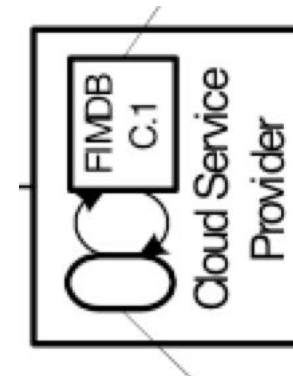2017

29

# Provided by the Cloud Service Provider

- File System
  - Managed services directory
  - OS binaries statically compiled for individual platforms

- Database
  - In-memory database landscape
  - Stored procedures and database algorithms
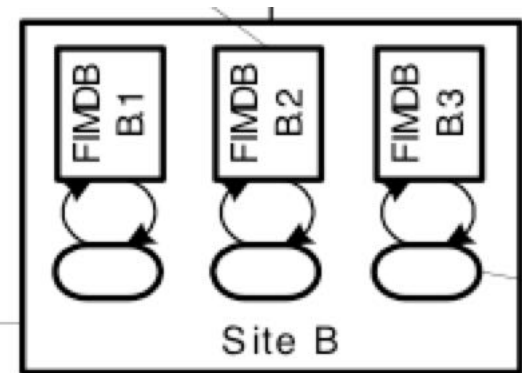  - Master application data

# Setup of a New Client

1. Establish site-to-site VPN connection b/w site and cloud service provider

2. Mount remote services directory

3. Install and configure local IMDB instance from services directory

4. Subscribe to and configure selected managed service

# Data Partitioning

- Supports parallel query execution
- Protects sensitive data
- Brings algorithms to data

**Details for Table**

| Parts | Columns | | |
| --- | --- | --- | --- |
| Host:Port/Partition ∧ | | Record Count | Total Size (KB) |
| ▼node-01:30203 | | | |
| 16 | | 85,286 | 2,675 |
| ▼node-02:30203 | | | |
| 15 | | 128,417 | 15,577 |
| ▼node-09:30203 | | | |
| 2 | | 78,873 | 2,489 |
| ▼node-10:30203 | | | |
| 8 | | 184,010 | 5,436 |
| ▼node-11:30203 | | | |
| 21 | | 112,729 | 3,252 |
| ▼node-14:30203 | | | |
| 13 | | 43,296 | 1,765 |
| ▼node-15:30203 | | | |
| 5 | | 93,507 | 3,075 |
| ▼node-17:30203 | | | |
| 7 | | 175,184 | 5,347 |
| ▼node-18:30203 | | | |
| 10 | | 270,924 | 28,734 |

**Use Case Oncology
Data Management**

Data Management for
Digital Health, Summer
2017

# Scheduling and Execution of GDPPs

1. Trigger task execution

**Webservice**

## Tasks

| ID | Pipeline | Params |
|----|----------|-----------|
| 12 | BWA | xyz.fastq |
| 13 | Stanford | A_1.fastq |
| 14 | Bowtie | xyz.fastq |

2. Schedule subtasks

**Scheduler**

### In-Memory Database

3. Execute subtasks

**Worker**

. . .

**Worker**

## Subtasks

| Task | ID | Job | Status | Params |
|------|----|--------|--------|-----------|
| 12 | 97 | Split | done | xyz.fastq |
| 12 | 98 | Import | todo | abc.vcf |
| 12 | 98 | Import | done | abc.vcf |

33

# Software Components and Communication

# IMDB Structure
# TASKS

**Table Name:**

**TASKS**

Columns | Indexes | Further Properties | Runtime Information

| | Name | SQL Data Type | Dimens | Column Store Data Type | Key | Not Null | Default |
|---|---|---|---|---|---|---|---|
| 1 | ID | INTEGER | | INT | | X | |
| 2 | STATUS | BIGINT | | FIXED | | | |
| 3 | PIPELINE_ID | INTEGER | | INT | | X | |
| 4 | PARAMETERS | VARCHAR | 5000 | STRING | | | |
| 5 | FASTQ_READCOUNT | BIGINT | | FIXED | | | 0 |
| 6 | CREATED_AT | TIMESTAMP | | LONGDATE | | | |
| 7 | USER | INTEGER | | INT | | | -1 |

▼ ⬛ WORKER
  ▦ JOBSTATISTICS
  ▦ NODE_GROUPS
  ▦ PIPELINES
  ▦ SESSIONS
  ▦ SUBTASKS
  ▦ TASKS

| 1,925 | 0 | 11 | {"file":{"name":"smallexample.fastq___2894","pretty_name":"smallexample.fastq","type":"use... | 0 | Oct 27, 2015 5:19:03.769 PM |
| 1,924 | 2 | 73 | {"file":{"name":"patient_NA12878_run_SRR000960.fastq___333970","type":"user","user_id":... | 333,970 | Oct 27, 2015 3:24:10.865 PM |
| 1,921 | 2 | 68 | {"file":{"name":"patient_NA12878_run_SRR000960.fastq___333970","type":"user","user_id":... | 333,970 | Oct 27, 2015 8:12:58.572 AM |
| 1,920 | 2 | 69 | {"file":{"name":"patient_NA12878_run_SRR000960.fastq___333970","type":"user","user_id":... | 333,970 | Oct 27, 2015 8:12:48.055 AM |
| 1,899 | 2 | 69 | {"file":{"name":"ERR005584.filt.fastq___65316","type":"user","user_id":3},"read_count":65316} | 65,316 | Oct 23, 2015 10:18:36.615 PM |
| 1,895 | 2 | 69 | {"file":{"name":"417kb.fastq___2263","type":"user","user_id":1},"read_count":2263} | 2,263 | Oct 23, 2015 10:01:09.444 PM |
| 1,894 | 3 | 68 | {"file":{"name":"417kb.fastq___2263","type":"user","user_id":1},"read_count":2263} | 2,263 | Oct 23, 2015 9:51:25.764 PM |

88 3 3 3 3 1 1

**Use Case Oncology Data Management**

Data Management for Digital Health, Summer 2017

35

# IMDB Structure
# SUBTASKS

# Runtime Layer
# Worker

- Workers execute jobs one by one
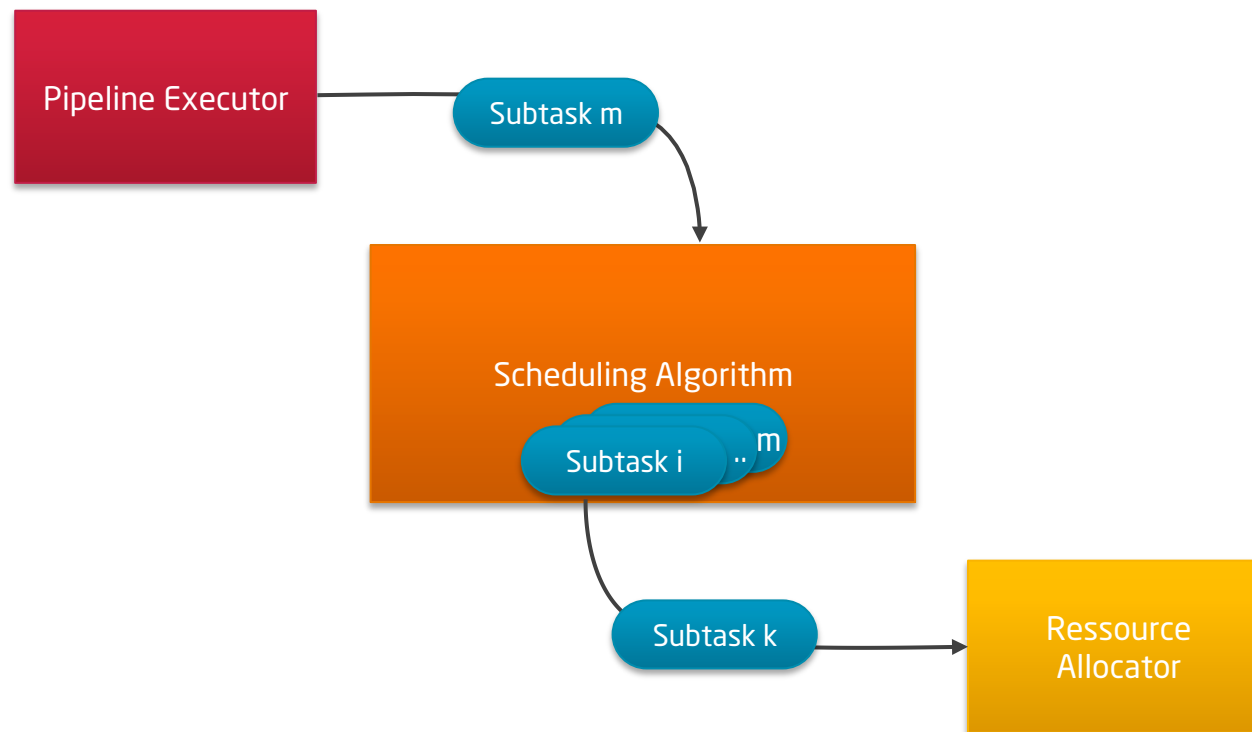- Subtask execution status in IMDB:
  - Ready (0),
  - In Progress (1),
  - Done (2), or
  - Erroneous (3).

- Jobs implemented as Python modules/classes
  - Can contain arbitrary code
  - Have access to IMDB
  - Can read/write to shared working directory

**Worker**

**IMDB**

**Node**
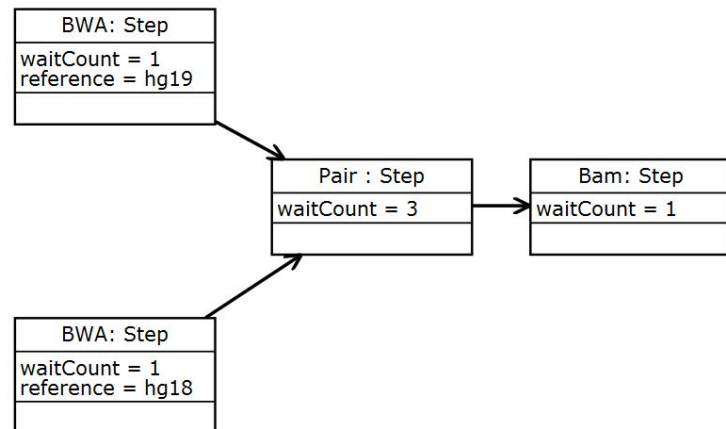
# Coordination Layer
## Scheduler

# Scheduling
# Pipeline Executer

```xml
<xpdl:Activity CompletionQuantity="1" Id="newpkg1_wp1_act2" Name="BWA">
    <xpdl:Implementation>
        <xpdl:No/>
    </xpdl:Implementation>
    <xpdl:Performers>
        <xpdl:Performer>newpkg1_wp1_par1</xpdl:Performer>
    </xpdl:Performers>
    <xpdl:NodeGraphicsInfos>
        <xpdl:NodeGraphicsInfo BorderColor="#000000" FillColor="#99FF99
            <xpdl:Coordinates XCoordinate="239.0" YCoordinate="219.0"/>
        </xpdl:NodeGraphicsInfo>
    </xpdl:NodeGraphicsInfos>
</xpdl:Activity>
```



| BWA: Step |
|---|
| waitCount = 1 |
| reference = hg19 |

| Pair : Step |
|---|
| waitCount = 3 |

| Bam: Step |
|---|
| waitCount = 1 |

| BWA: Step |
|---|
| waitCount = 1 |
| reference = hg18 |

# Scheduling Algorithms

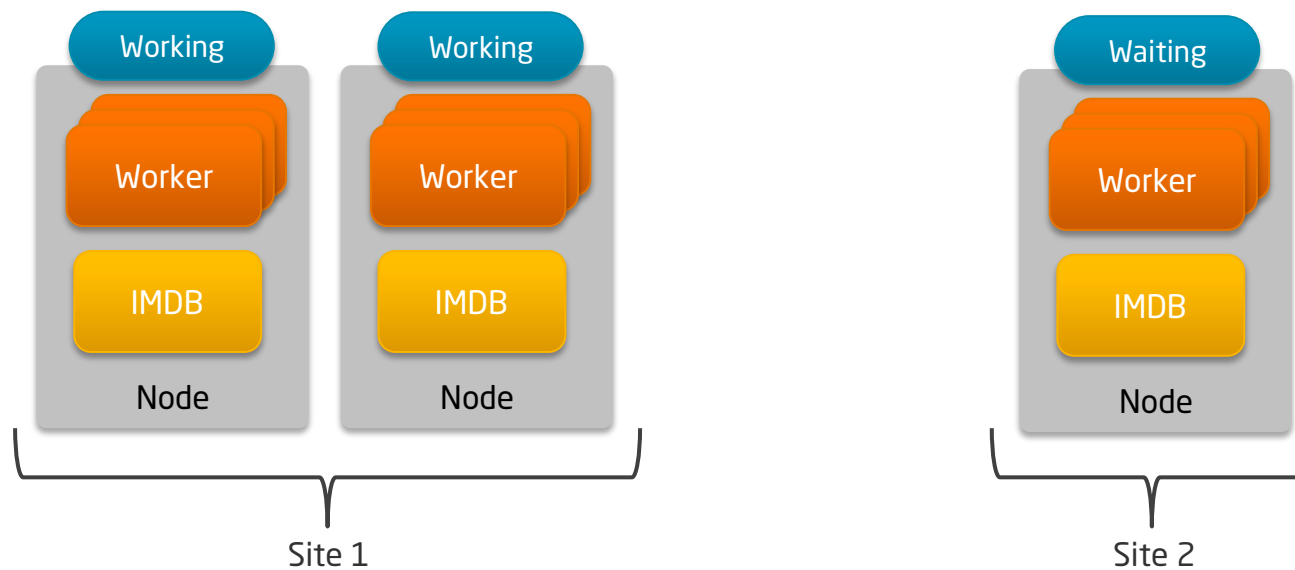- Scheduling algorithms are plug-in software modules

  - "User-/Group-based" to let users execute their tasks on their local site only

  - "Priority First" to prefer important users

  - "High Throughput", i.e. "shortest task first" to deal with high load

- Scheduling algorithms can also be composed hierarchically

# Scheduling
# Resource Allocator

- Maintains lists of running and idle nodes

- Idle worker requests new sub task for its assigned groups

- If there is no matching sub task, it sleeps until a new sub task gets ready
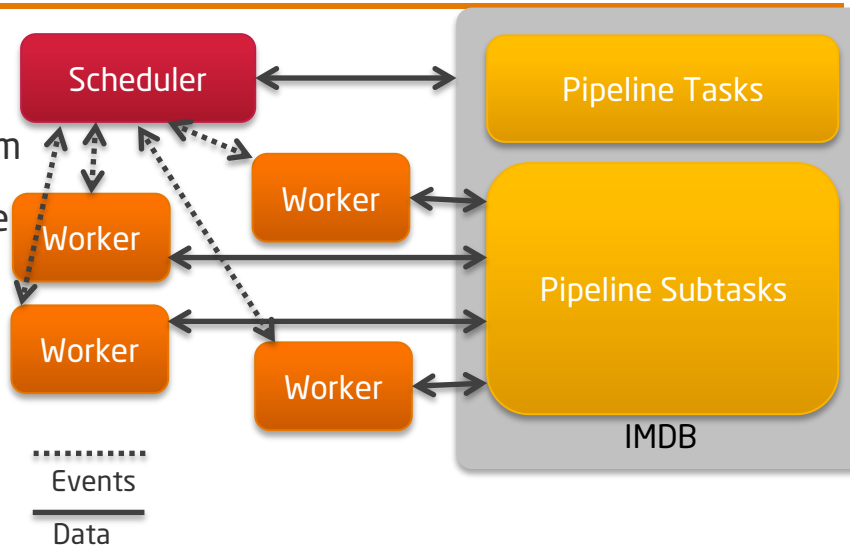
# Recoverability

- All execution data is stored in IMDB

- Temporary files on a shared file system

- In case of any failure, the system-wide state can be restored

Scheduler

Worker

Worker

Worker

Worker

Pipeline Tasks

Pipeline Subtasks

IMDB

......... Events

Data

```
TypeError:  NoneType  object is unsubscriptable

2015-11-04 18:01:30 INFO       [ContinuingCoordinator] will start task with ID 1969
2015-11-04 18:01:30 INFO       [ContinuingCoordinator] Will continue old but unfinished task 1969 with 52 already done subtasks.
2015-11-04 18:01:31 ERROR      [ContinuingCoordinator] Traceback (most recent call last):
```

# FIMDB Comparison

- FIMDB provides (smaller) algorithms to (larger) data
- Forms a single virtual database across sites and locations
- Master data managed by service provider whilst sensitive data resides locally

| Pros | Cons |
|---|---|
| Single database license | Complex operation |
| Easy to consume services | Single setup required |
| Query propagation by IMDB | |