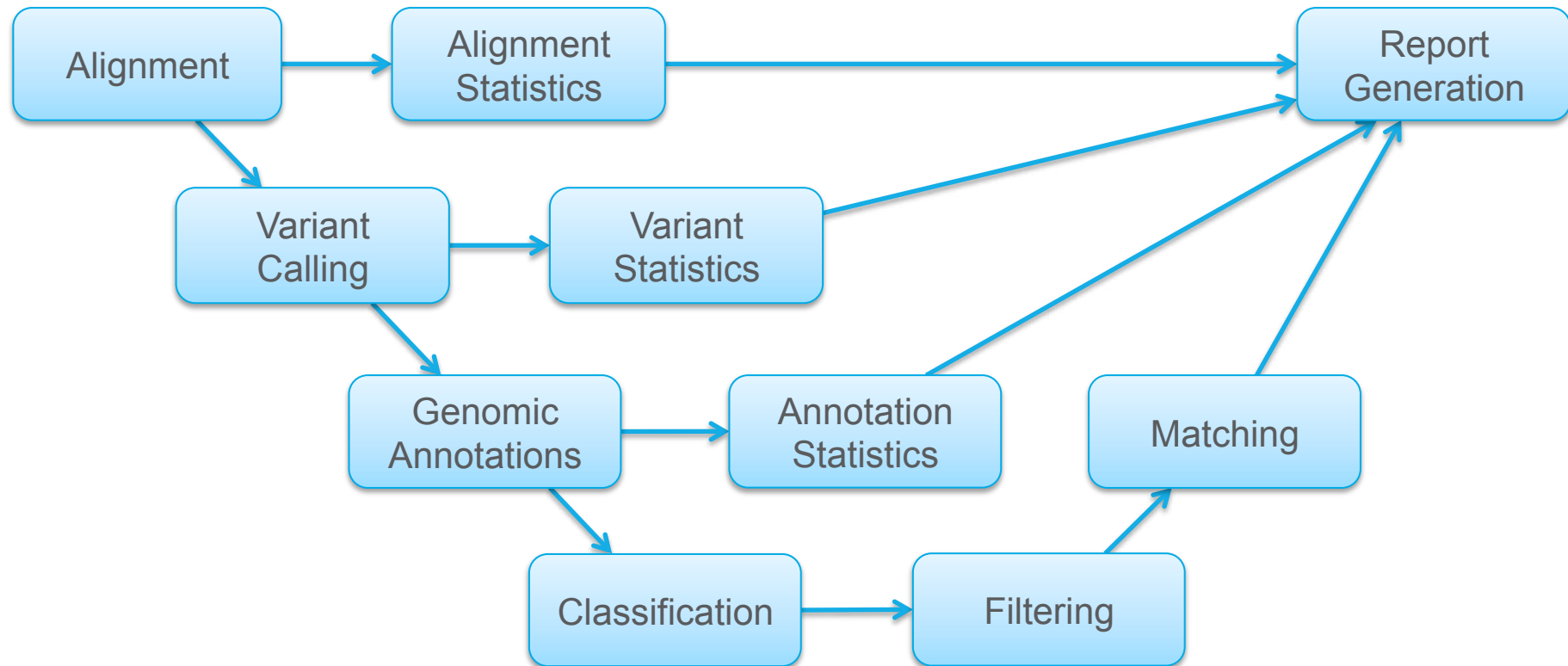# Data-Flow Programming

Motivation

Related Work

CloudKeeper

# Example: Genome-Analysis Workflows

# Related Work

## Academic Workflow Systems
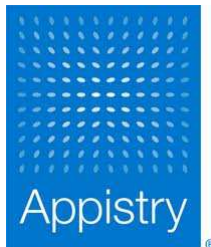
Taverna

UNICORE

Pegasus

Appistry

nextflow

Galaxy

## Business Process Execution Language

ORACLE
BPEL PROCESS MANAGER
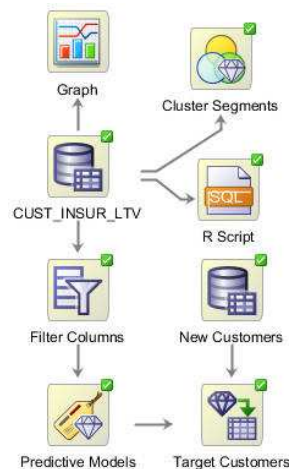
## Batch Processing

- Java EE 7 / JSR 352

- spring Batch, XD

- Apache Spark

## Domain-Specific Workflow Tools

- Oracle Data Miner

Graph

Cluster Segments

CUST_INSUR_LTV

R Script

Filter Columns

New Customers

Predictive Models

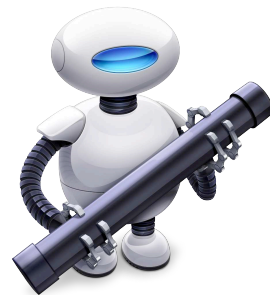Target Customers

# Three Select Aspects for Taxonomy

Embedded as library                    Stand-alone application

Textual programming language

Graphical user interface

Binary artifact repository
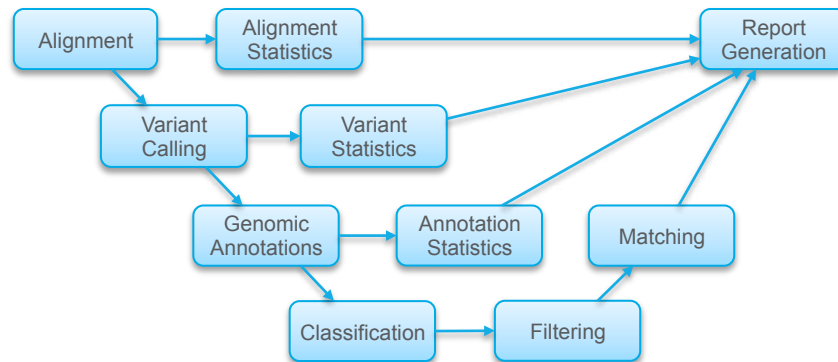
Social-network-like web site

Software Engineers                    Users: Geneticists, etc.

# CloudKeeper
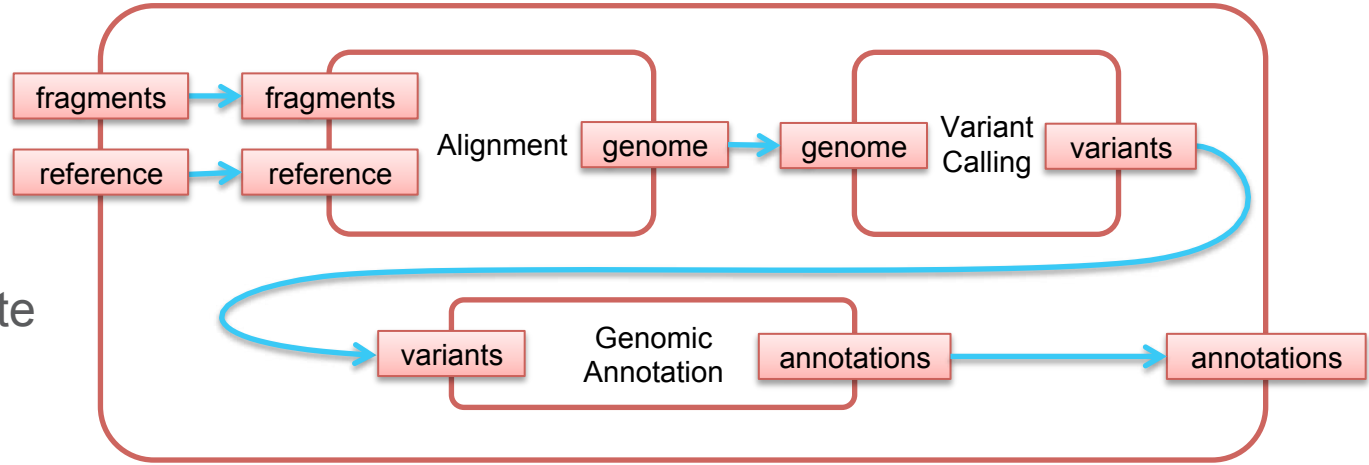
## Lifecode's Data-Flow Language and Runtime System

**1** Modular data-flow programming "in the large"

**2** Scalable: write once, run everywhere

**3** Satisfy needs of software engineers and expert users

**4** Lightweight and easy to deploy

# Modular Data-Flow Programming "in the Large"

## Modules

- Functional units

- In- and out-ports

- Simple or composite



## Programming "in the Large"

- Instantiate modules and define connections between ports

- Control flow via special composite modules (e.g., loop modules)

# 2 Scalable: Write Once, Run Everywhere

## Scale Horizontally

• Runs within a single Java Virtual Machine on a laptop

• Same code also runs on a cluster or in the cloud

**DRMAA**
Distributed Resource Management
Application API — www.drmaa.org

**GRID ENGINE**

**HTCondor**
High Throughput Computing

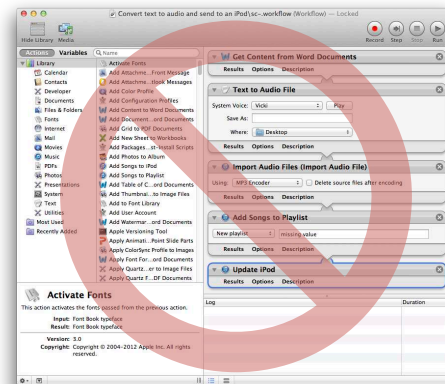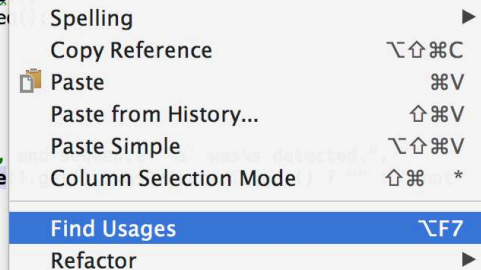• Transparent data transfer using main memory, file system, S3, etc.

• Checkpointing

# 3 Satisfy Needs of Software Engineers and Expert Users

## Workflow System with Developers in Mind

- Internal DSL on top of Java

  - Works with all IDEs

  - Easy refactoring

- Statically typed and verified



```
public abstract class ReportModule extends SimpleModule<ReportModule> {
    public abstract InPort<Double> avgLineLength();
    public abstract InPort<String> subsequence();
    public abstract InPort<Boolean> wasDetected();
    public abstract OutPort<String> report();

    @Override
    public void run() {
        report().set(String.format(
            "Report: Avg. read length is %.2f,
            avgLineLength().get(), subsequence
        ));
    }
}
```

| | |
|---|---|
| Spelling | ▶ |
| Copy Reference | ⌥⇧⌘C |
| Paste | ⌘V |
| Paste from History... | ⇧⌘V |
| Paste Simple | ⌥⇧⌘V |
| Column Selection Mode | ⇧⌘* |
| **Find Usages** | ⌥F7 |
| Refactor | ▶ |

```
17
18   public class ITGenomeAnalysis {
19   ▶ Run 'ITGenomeAnalysis'        ^⇧F10
20   🐞 Debug 'ITGenomeAnalysis'      ^⇧F9
21   ▒ Run 'ITGenomeAnalysis' with Coverage
22
23       public void setup() {
24           cloudKeeper = new SingleVMCloudKeeper.
25           cloudKeeperEnvironment = cloudKeeper.n
26       }
27
```

# 4 Lightweight and Easy to Deploy

## CloudKeeper Embedded

- Library, not framework

- Debug in single Java Virtual Machine

- High-level alternative to `ExecutorService`, actors, futures, etc.

## Package Management

- Execute directly from artifact repository

  – Similar to Groovy's Grape
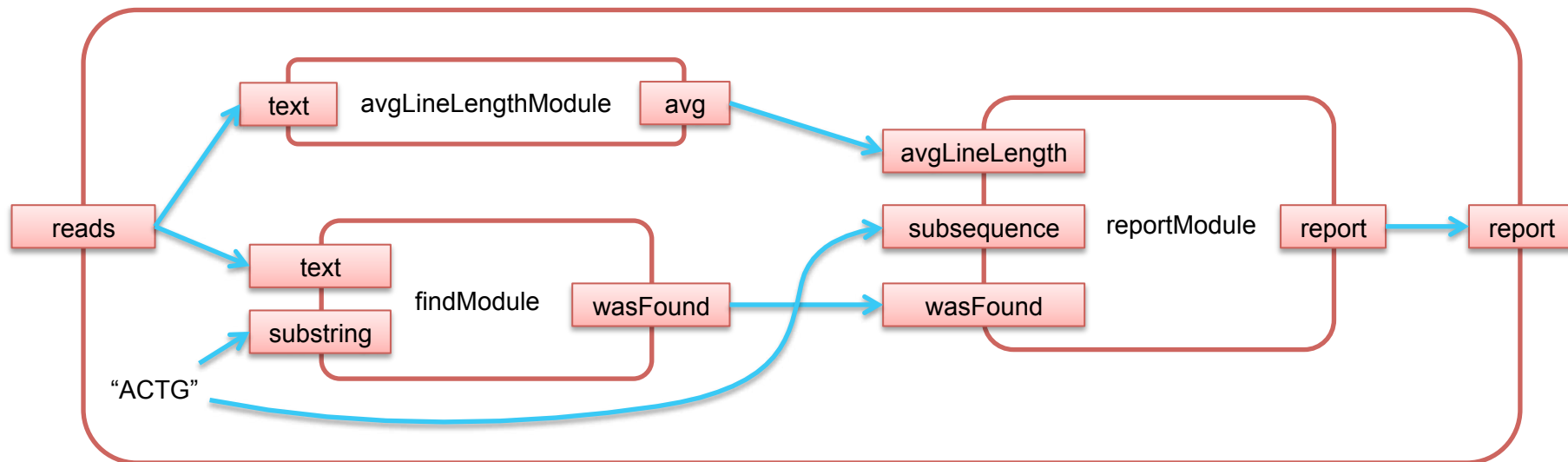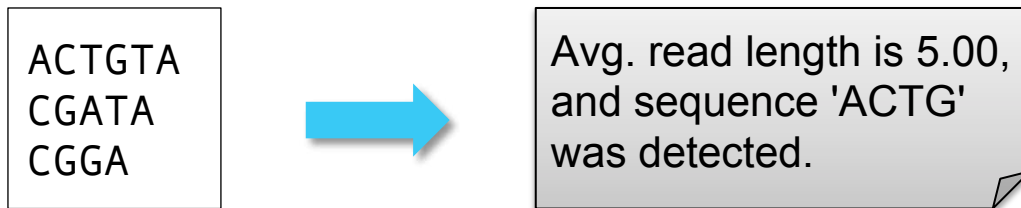
eclipse Aether

Nexus

artifactory

# Example

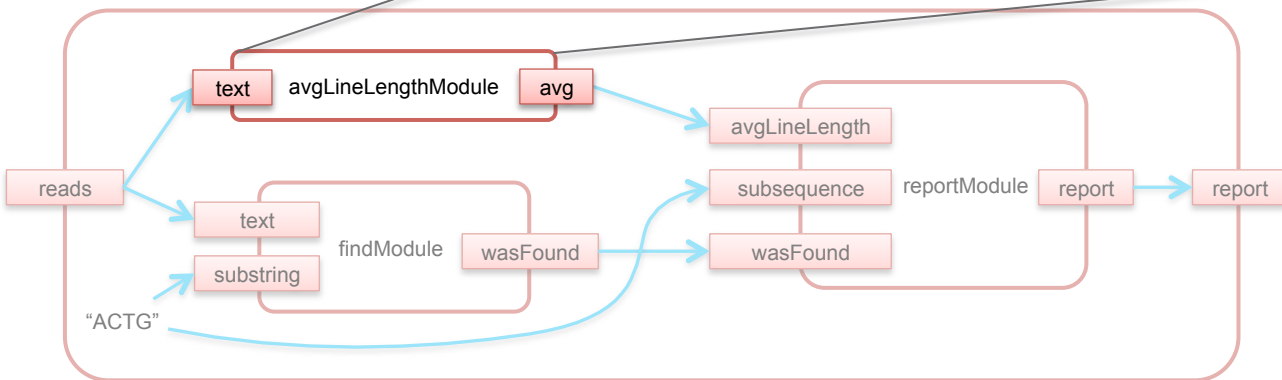Simple Modules

Composite Modules

# Example

# Simple Modules

## Internal DSL

- Module ~ Java class

- Port ~ method

- Metadata using annotations

```java
@SimpleModulePlugin("Computes the average line length in a text")
public abstract class AvgLineLengthModule
        extends SimpleModule<AvgLineLengthModule> {
    public abstract InPort<String> text();
    public abstract OutPort<Double> avg();

    @Override
    public void run() throws IOException {
        String text = text().get();
        double avg;
        // ...
        avg().set(avg);
    }
}
```

# Composite Modules

```java
@CompositeModulePlugin("Analyzes String consisting of DNA fragments")
public abstract class GenomeAnalysisModule
        extends CompositeModule<GenomeAnalysisModule> {
    public abstract InPort<String> reads();
    public abstract OutPort<String> report();

    InputModule<String> sequence = value("ACTG");
    AvgLineLengthModule avgLineLengthModule = child(AvgLineLengthModule.class)
        .text().from(reads());
    FindModule findModule = child(FindModule.class)
        .text().from(reads())
        .substring().from(sequence);
    ReportModule reportModule = child(ReportModule.class)
        // ...

    { report().from(reportModule.report()); }
}
```