

# Algorithmic Infrastructure and Performance of Finding Optimal Conditional Sparse $l_p$ -norm Regression

John Hainline

# Outline

- Problem Overview
- Formulation
- Algorithms
- Implementation
- Architecture
- Evaluation
- Conclusion

# Problem Overview

- Goals
  - Provide scalable and robust implementation of new algorithm in machine learning
  - Specifically the "Conditional Sparse  $L_p$ -norm Regression Problem" (explained next)
  - Focused on engineering a performant parallelized cloud solution
  - Created database to provide complete information recovery

# Problem Overview

We are trying to implement an algorithm that solves a Conditional Sparse  $l_p$ -norm Regression problem

- Given  $X \in \mathbb{B}^{n \times p} = (\vec{x}_1, \vec{x}_1, \dots, \vec{x}_n)^T$   
 $Y \in \mathbb{R}^{n \times m} = (\vec{y}_1, \vec{y}_2, \dots, \vec{y}_n)^T$   
 $Z \in \mathbb{R}^n$
- Find conditional  $\mathbf{c}$  (our k-DNF) such that  $\mathbf{c}(\mathbf{X})$  creates a subpopulation where the regression line  $\langle \vec{a}, Y \rangle$  predicts  $\mathbf{Z}$ 
  - $\mathbf{a}$  is sparse, i.e. uses few components of  $\mathbf{Y}$  (e.g. 2)
- Using  $l_2$ -norm (i.e. least squared error)

# Problem Overview

- **k Disjunctive Normal Form (k-DNF)**
  - A formula is a DNF iff it is a disjunction ( $\vee$ ) of conjunctions ( $\wedge$ ) of literals
  - A **k**-DNF means all conjunctions contain **k** literals
  - Valid 2-DNFs
    - $(x_1 \wedge x_2)$
    - $(x_8 \wedge \neg x_2) \vee (\neg x_2 \wedge x_5) \vee (\neg x_1 \wedge x_2) \vee (x_3 \wedge x_4)$
  - Valid 3-DNFs
    - $(x_1 \wedge x_2 \wedge x_3)$
    - $(x_1 \wedge \neg x_2 \wedge \neg x_3) \vee (\neg x_1 \wedge \neg x_2 \wedge x_5) \vee (\neg x_2 \wedge x_3 \wedge x_4)$
  - Using 2-DNFs although system could easily be extended for arbitrary k-DNFs

# "Red-Blue Set Cover" Formulation

- Working with real-valued variables, and using a greedy strategy
  - "Sets" are all possible terms  $(\mathbf{x}_i \wedge \mathbf{x}_j)$  of which there are:  $4 \times \binom{p}{2}$
  - i.e. each term denotes a set of some points of  $\mathbf{Y}$
  - Originally couched in terms of red or blue categories for points
    - Required to cover some number of blue points
    - Number of red points covered is the "cost"
  - Construct  $\mathbf{c}(\mathbf{X})$  (our 2-DNF) such that we:
    - Minimize "redness" of points (i.e. distance from regression line  $\langle \vec{a}, Y \rangle$ )
    - Maximize number of points

# Optimizations

- Sub-sampling
  - Reduce size of  $\mathbf{Y}$  using sub-sampling ( $\mathbf{n} \rightarrow \mathbf{r}$ )
  - High probability of finding hidden k-DNFs
- Combinatorics
  - Must compute Set Cover  $\binom{r}{2} \times \binom{m}{2}$  times

# Algorithms

---

**Algorithm 1:** This algorithm breaks down the relevant search space of  $X \in \mathbb{B}^{n \times p}$  and  $Y \in \mathbb{R}^{n \times m}$ , giving SetCover the parameters it needs. Note that *terms* must include the  $\neg$  versions of our  $\binom{p}{2}$  possible combinations of  $X$ , which is why we have  $4 \cdot \binom{p}{2}$  combinations.

---

```
1 def ComputeBestKDNF( $X, Y, Z, r$ ):
2    $X, Y = \text{SubsamplePoints}(X, Y, r)$ 
3    $terms = (x_i \cap x_j)$  for all  $4 \cdot \binom{p}{2}$  combinations of  $X$ 
4    $best\_terms = \emptyset$ 
5    $best\_error = \infty$ 
6   for  $(i, j) \in \binom{r}{2}$  points do
7     for  $(k, l) \in \binom{m}{2}$  dimensions do
8        $t, e = \text{SetCover}(X, Y, Z, terms, i, j, k, l)$ 
9       if  $e < best\_error$  then
10          $best\_terms = t$ 
11          $best\_error = e$ 
12       end
13     end
14   end
15   return  $best\_terms, best\_error$ 
```

---



# Algorithms

---

**Algorithm 2:** The SetCover algorithm searches for the best  $k$ -DNF (lowest error in terms of "redness") for a given pair of data points  $i, j$  and dimensions  $k, l$  in  $Y \in \mathbb{R}^{n \times m}$ . The variable *terms* holds every *set* generated by our  $\binom{p}{2}$  terms  $(x_i \cap x_j)$  in  $X \in \mathbb{B}^{n \times p}$ .  $\mu$  dictates how large of a sub-population we are looking for.

---

```
1 def SetCover( $X, Y, Z, terms, i, j, k, l$ ):
2    $\vec{a}$  = CalculateRegressionLine( $\vec{y}[k]_i, \vec{y}[l]_j$ )
3    $r\_points$  = CalculateRednessForPoints( $Y, \vec{a}$ )
4    $r\_terms$  = CalculateRednessForTerms( $terms, r\_points$ )
5    $sort$  = SortByAverageRedness( $r\_terms$ )
6    $best\_terms = \emptyset$ 
7    $best\_error = \infty$ 
8    $cur\_terms = \emptyset$ 
9   while ( $cur\_terms.size < \mu \cdot n$ ) do
10      $next$  = NextUnusedTermInSort( $sort, cur\_terms$ )
11      $cur\_terms.add(next)$ 
12   end
13    $cur\_error$  = AverageRednessPerPoint( $cur\_terms$ )
14   if  $cur\_error < best\_error$  then
15      $best\_error = cur\_error$ 
16      $best\_terms = cur\_terms$ 
17   end
18   return  $best\_terms, best\_error$ 
```

---

# Implementation

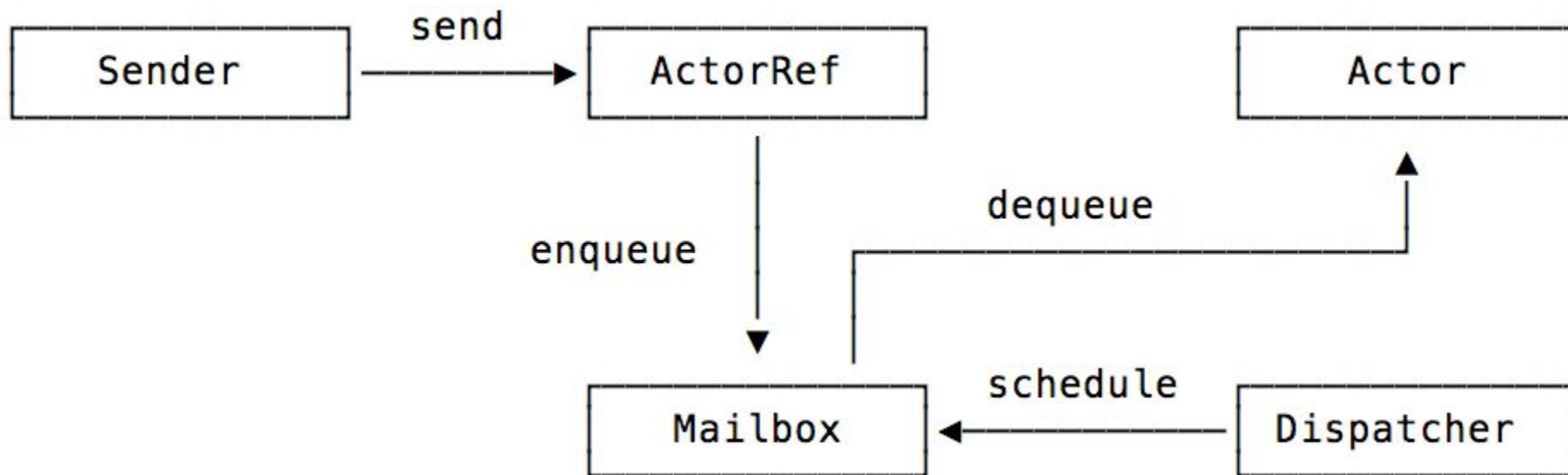
- Akka Library
  - Akka Actors
  - Akka Streams
  - Akka Http
  - Akka Cluster
  - Cluster Sharding
  - Distributed Data
  - Akka Persistence
  - Alpakka
  - Akka gRPC
  - Akka Management

# Implementation

- Akka Library
  - **Akka Actors**
  - **Akka Streams**
  - Akka Http
  - **Akka Cluster**
  - Cluster Sharding
  - Distributed Data
  - Akka Persistence
  - Alpakka
  - Akka gRPC
  - Akka Management

# Implementation

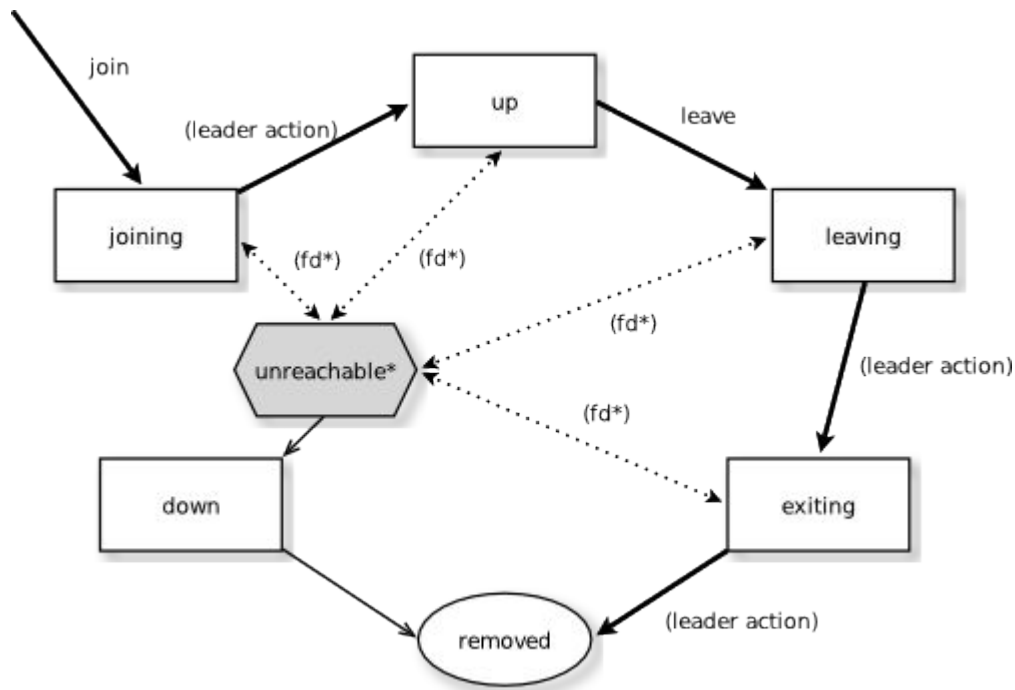
- Akka Actors



Source: <https://blog.codecentric.de/en/2015/08/introduction-to-akka-actors/>

# Implementation

- Akka Cluster
  - Member States
    - Joining
    - Up
    - Leaving
    - Exiting
    - Removed
    - Down



Source: <https://doc.akka.io/docs/akka/2.5/common/cluster.html>

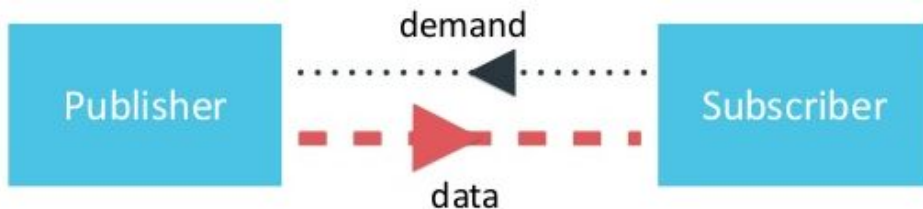
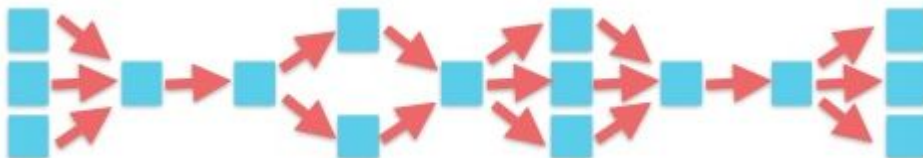
# Implementation

- Akka Streams

- Graph describes movement of data
- Starts at *Source(s)*
- Travels through arbitrarily complex graph of *Flows*
- Ends at *Sink(s)*

- Back-Pressure

- *Sink* demands data
- Demand moves backward to *Source*
- *Source* pushes data
- This manages stream speed and keeps nodes from desynchronizing



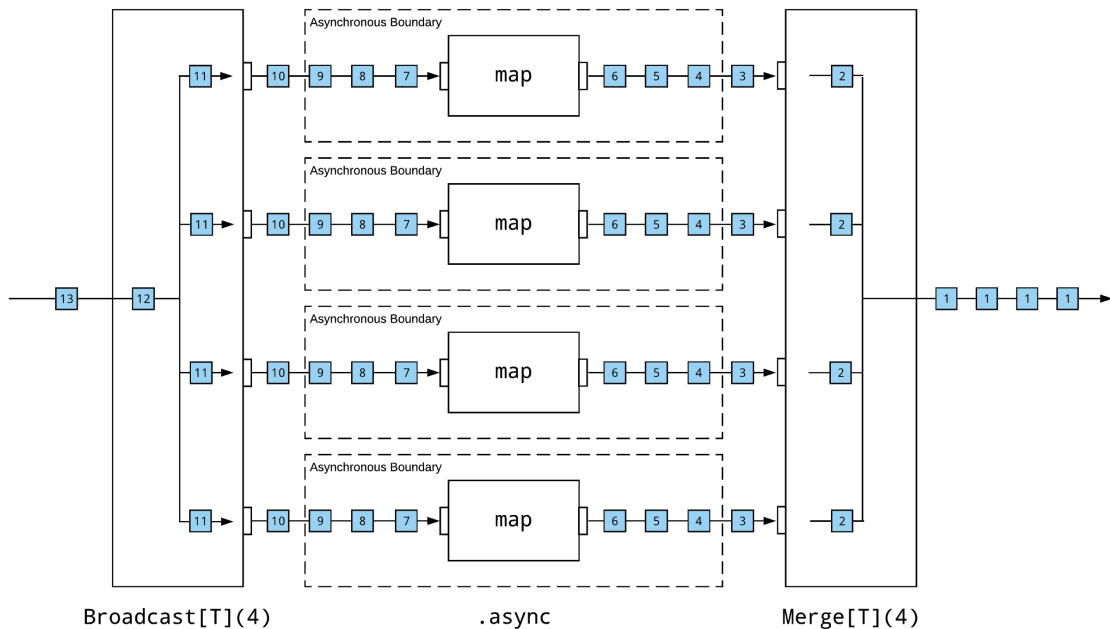
Source for images:

[https://www.slideshare.net/Typesafe\\_Inc/reactive-streams-100-and-why-you-should-care-webinar](https://www.slideshare.net/Typesafe_Inc/reactive-streams-100-and-why-you-should-care-webinar)

# Implementation

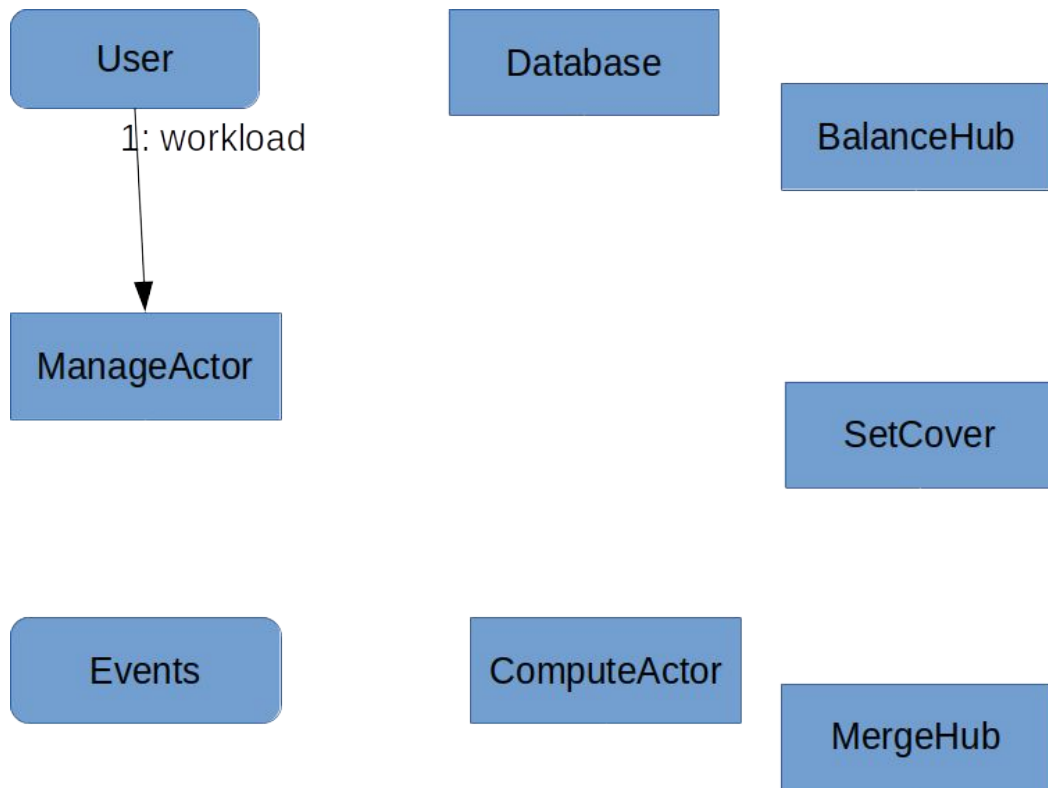
- Akka Streams

- *Source* sends **work** to *BalanceHub*
- *BalanceHub* broadcasts **work** to *SetCover* instances
- *SetCover* runs algorithm, sending on **results**
- *MergeHub* merges **results** to a single stream and saves



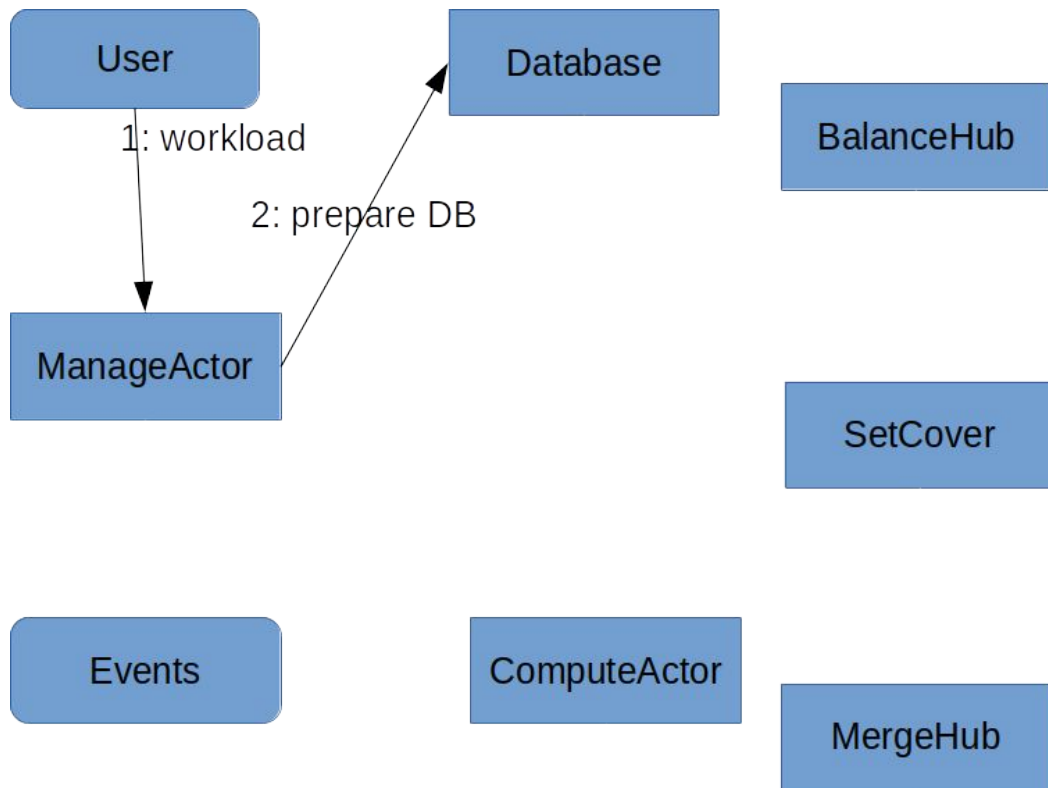
Source: <https://blog.colinbreck.com/partitioning-akka-streams-to-maximize-throughput/>

# Architecture

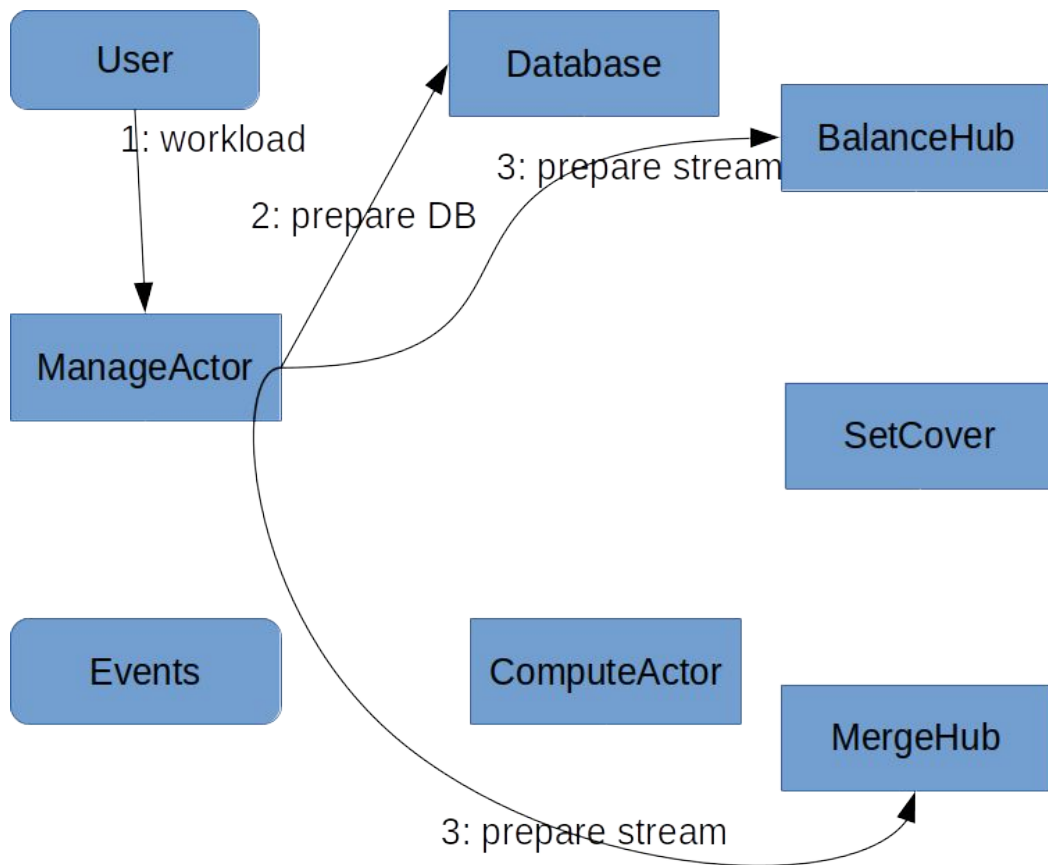




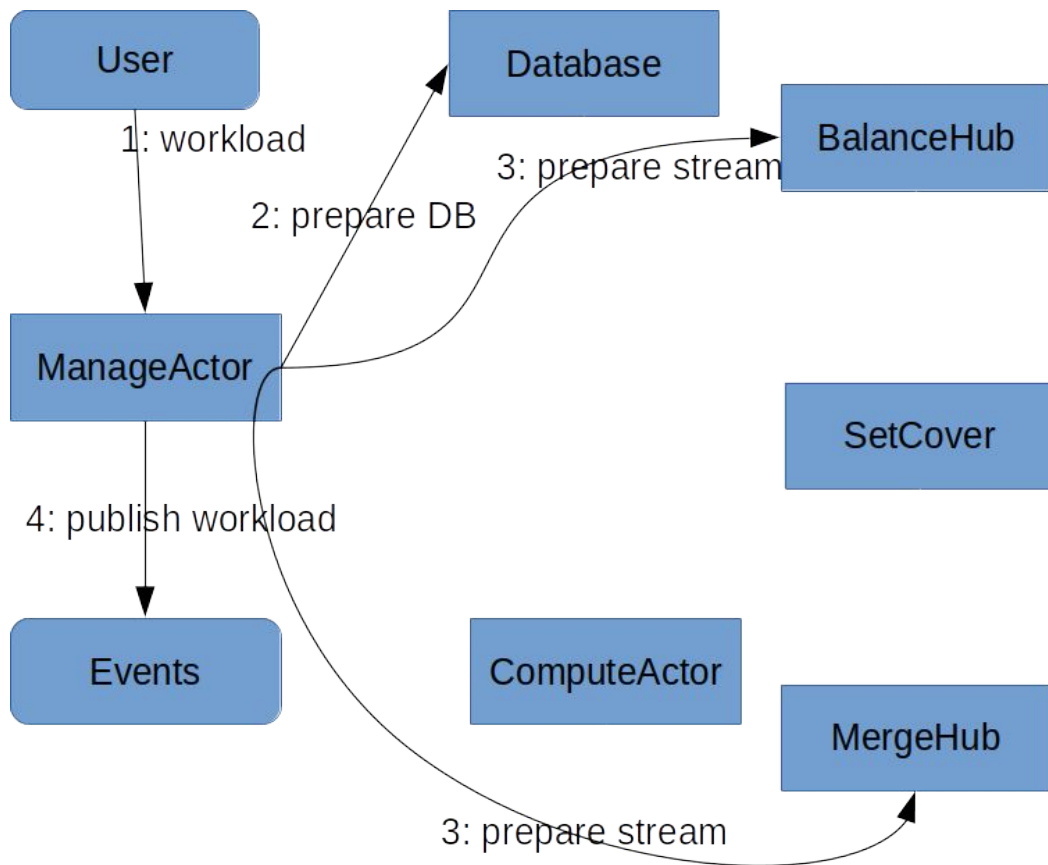
# Architecture



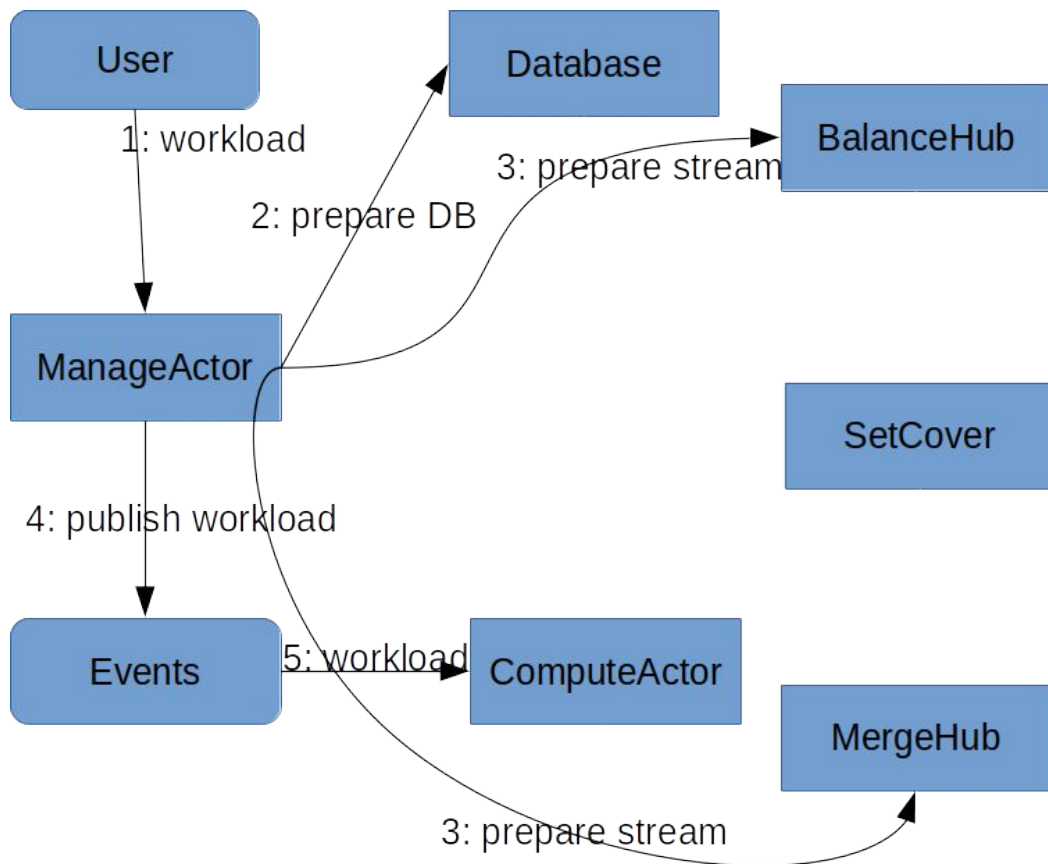
# Architecture



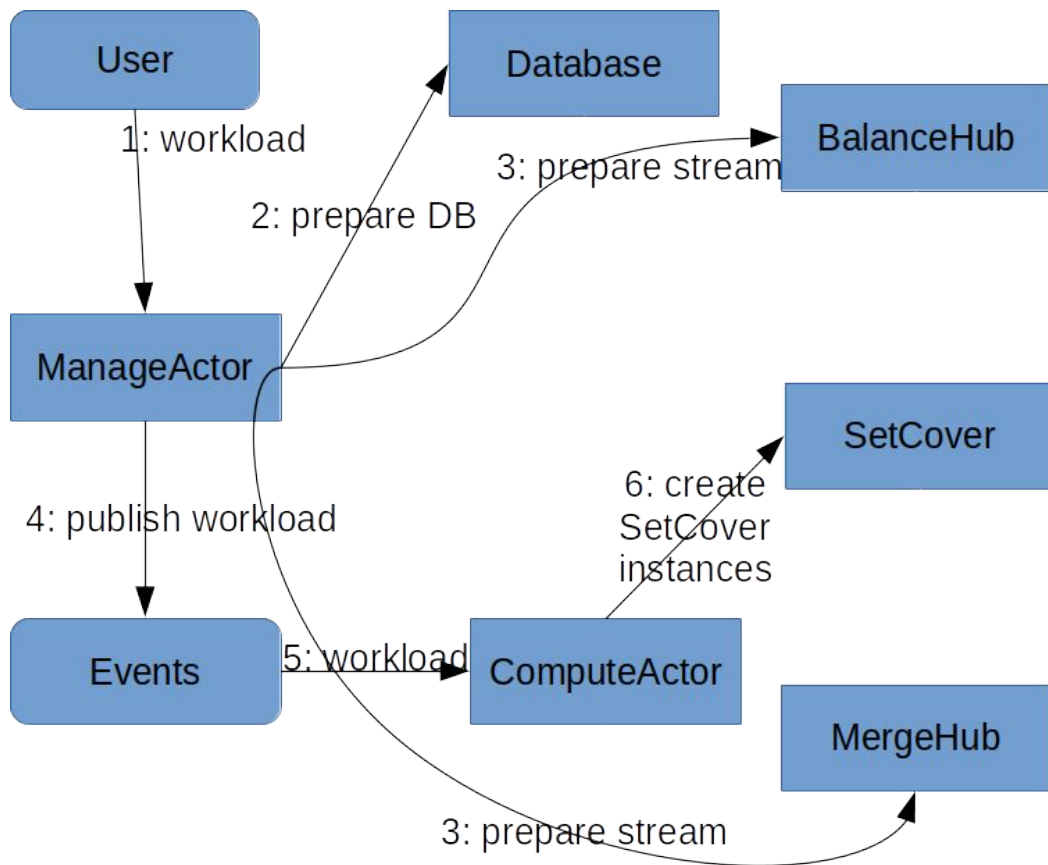
# Architecture



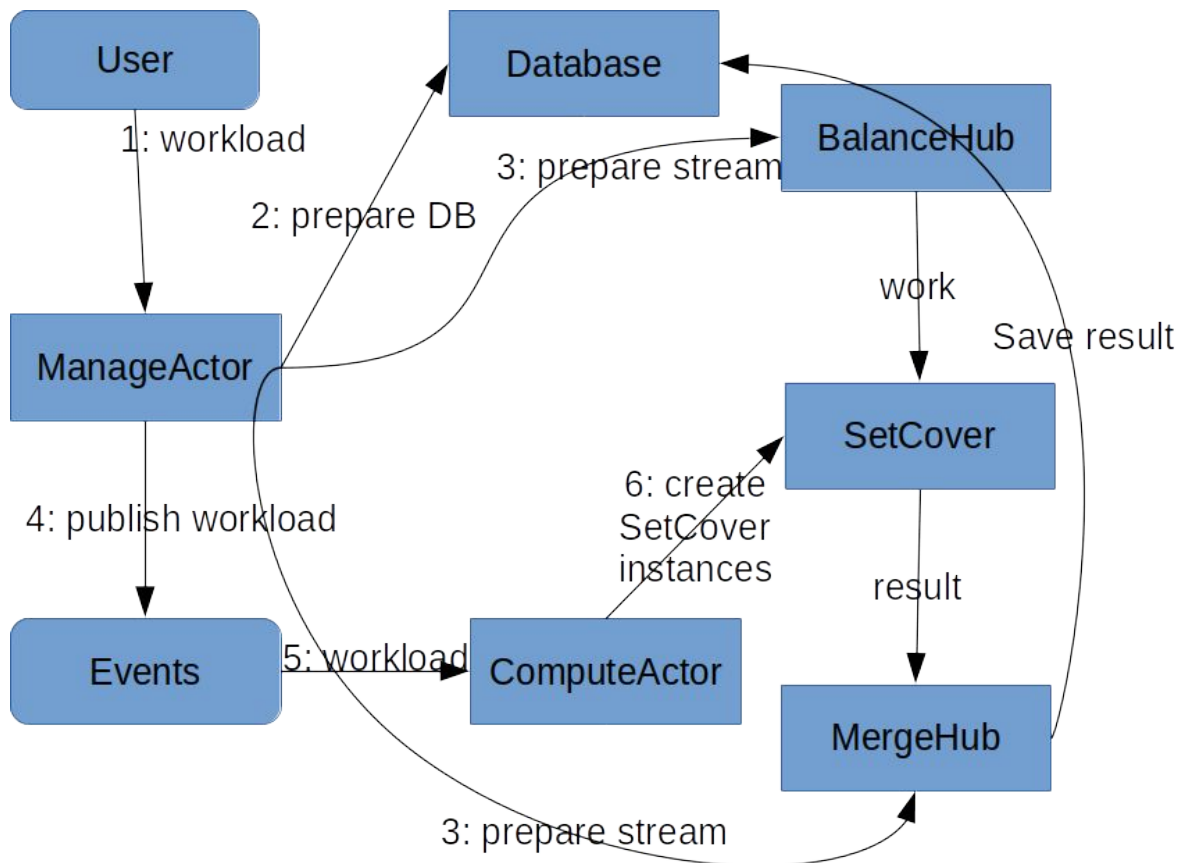
# Architecture



# Architecture



# Architecture



# Architecture

- Workload
  - name: data .sql filename
  - mu: % of data which k-DNF includes
  - optionalSubset: number of points to select from  $Y$
  - optionalRandomSeed: RNG seed for selecting  $Y$  points
- Work
  - index: absolute "position" in Co-Lex order
  - selectedDimensions: (k, l)
  - selectedRows: (i, j)
- Result
  - index: absolute "position"
  - selectedDimensions: (k, l)
  - selectedRows: (i, j)
  - coefficients: sparse regression line
  - kDNF:  $(x1 \cap \neg x2) \cup (x3 \cap x5)$

```
case class Workload(name: String, mu: Double,  
                    optionalSubset: Option[Int] = None,  
                    optionalRandomSeed: Option[Int] = None)  
  
case class Work(index: Int,  
                selectedDimensions: Array[Int],  
                selectedRows: Array[Int])  
  
case class Result(index: Int,  
                  dimensions: Array[Int],  
                  rows: Array[Int],  
                  coefficients: Array[Double],  
                  error: Option[Double],  
                  kDNF: Option[String])
```

# Evaluation

- Sub-sampling increased speed by multiple orders of magnitude
  - Uses a RNG seed so we can recover a failed workload
- BitSet implementation gave 10x speed increase
  - Went from using 30 CPUs to 100 CPUs
- Batching (in streams) gave 10x speed increase
  - Went from using 100 CPUs to 200+ CPUs
- Expectation
  - Likely scales to 300+ CPUs
  - Not tested due to cost on GKE



# Conclusion

- System dynamically links SetCover Flows into Akka Stream
  - Allows scaling across an entire cluster
- Using Sub-sampling, Batching, and BitSet for speed
- MySQL Database allows recovery on system failure
- Configured to run on Google Kubernetes Engine (or a local machine)
- Successfully ran
  - $n = 10,000$  points, sub-sampled to  $s = 500$
  - SetCover ran 27,702,500 times over 5 days
  - Used 208 CPUs across 26 VMs

# Conclusion

Questions?