

Building Streaming Recommendation Engines on Spark

Rui Vieira

rcardoso@redhat.com

Overview

- Collaborative Filtering
 - Batch Alternating Least Squares (ALS)
 - Streaming ALS
- Apache Spark
 - Distributed Streaming ALS
- OpenShift deployment

Collaborative Filtering

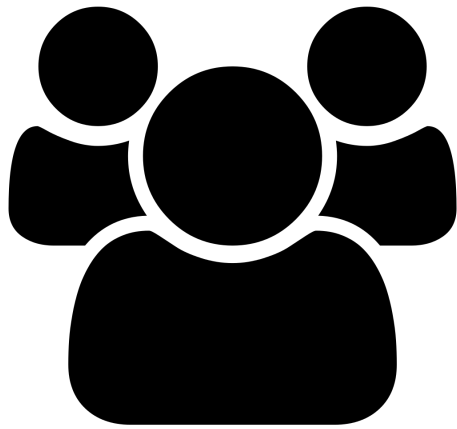
- Users, products and ratings

$(\text{user}, \text{product}) \rightarrow \text{rating}$

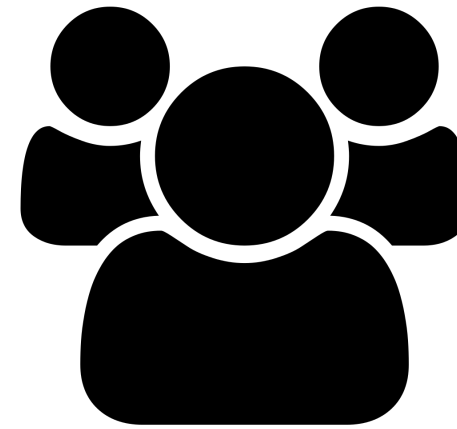
- Collaborative
- “Filtering”

Collaborative Filtering

A



B



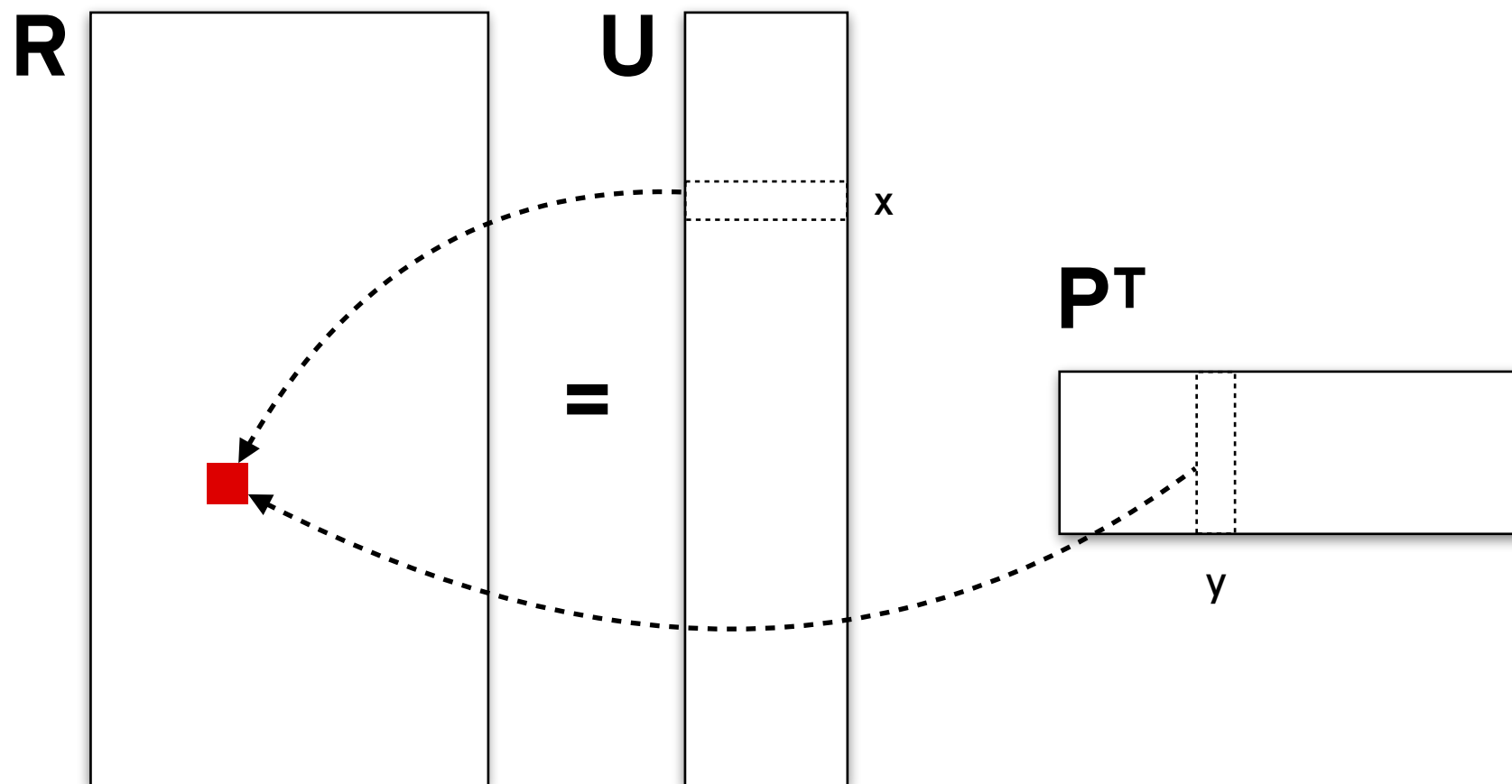
Collaborative Filtering



Alternating Least Squares

$$R = \begin{array}{ccccc} & \text{user 1} & \text{user 2} & \text{user 3} & \dots & \text{user N} & \\ \left[\begin{array}{c} 1 \\ ? \\ 5 \\ \vdots \\ 2 \end{array} \right. & \begin{array}{c} 4.5 \\ 3 \\ 3 \\ \vdots \\ 4 \end{array} & \begin{array}{c} ? \\ 3 \\ ? \\ \vdots \\ 1 \end{array} & \begin{array}{c} \dots \\ \dots \\ \dots \\ \ddots \\ \dots \end{array} & \begin{array}{c} 3 \\ 4 \\ ? \\ \vdots \\ ? \end{array} & \begin{array}{c} \text{product 1} \\ \text{product 2} \\ \text{product 3} \\ \vdots \\ \text{product M} \end{array} \end{array}$$

Alternating Least Squares



$$\hat{r}_{x,y} = U_x P_y^T$$

Batch ALS

$$\text{loss} = \sum_{x,y} \left(\underbrace{r_{x,y} - \hat{r}_{x,y}}_{\epsilon_{x,y}} \right)^2 + \lambda_x \sum_x ||\mathbf{U}_x||^2 + \lambda_y \sum_y ||\mathbf{P}_y||^2$$

Batch ALS

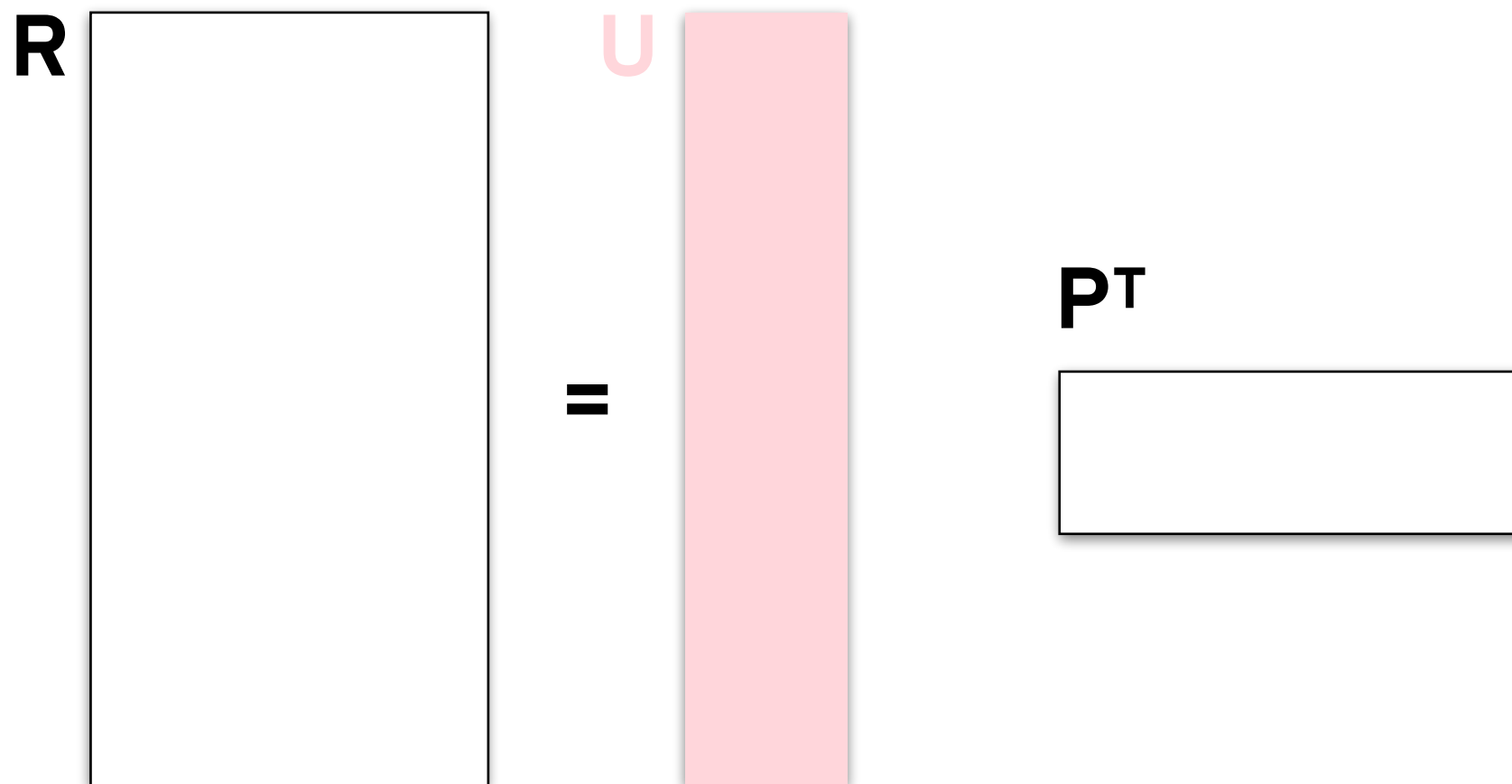
$$\text{loss} = \sum_{x,y} \left(\underbrace{r_{x,y} - \hat{r}_{x,y}}_{\epsilon_{x,y}} \right)^2 + \lambda_x \sum_x \|U_x\|^2 + \lambda_y \sum_y \|P_y\|^2$$

(minimize)

$$\frac{\partial \text{loss}}{\partial U_x} = 0,$$

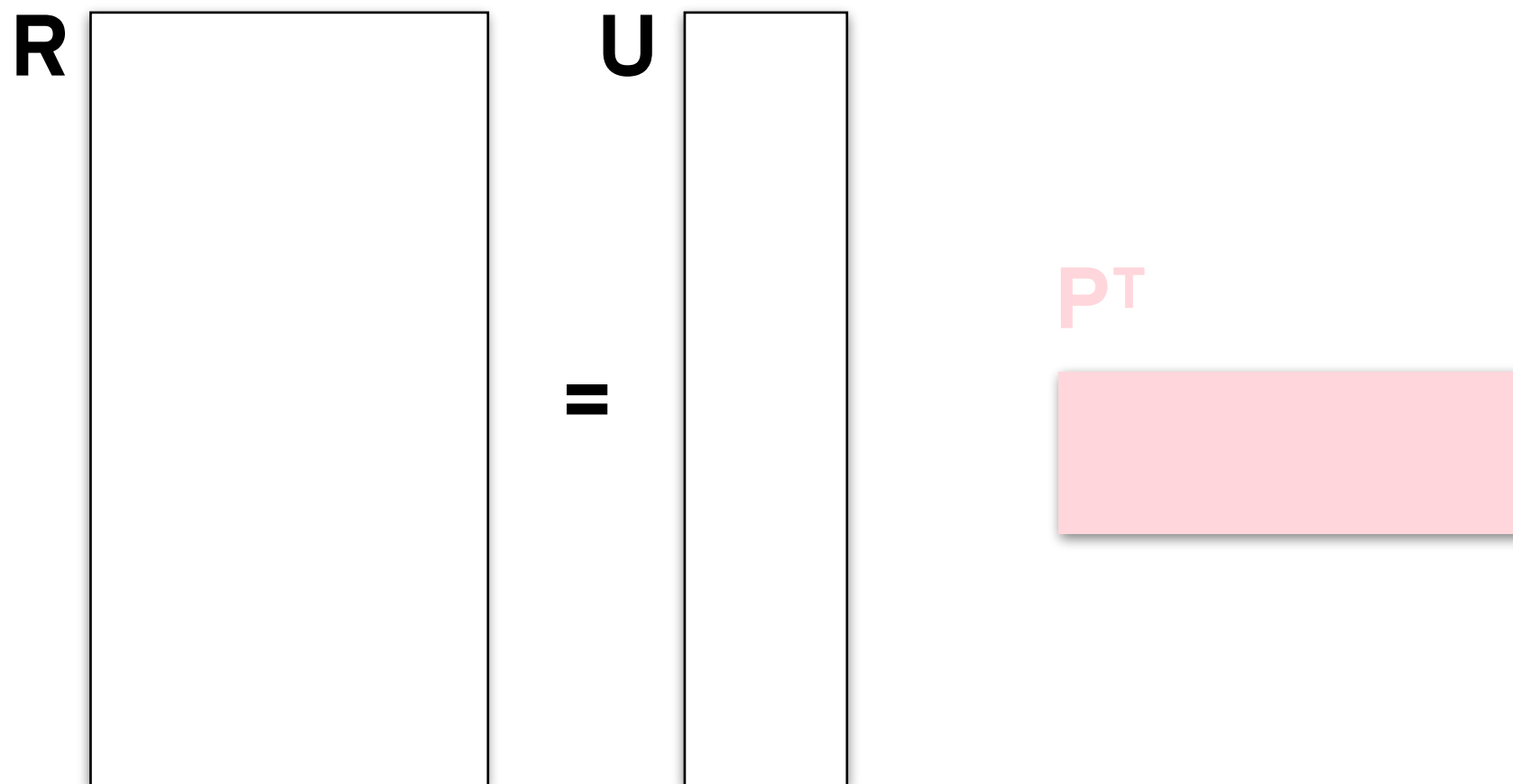
$$\frac{\partial \text{loss}}{\partial P_y} = 0$$

Alternating Least Squares



$$P_y = r_y X (X^T X + \lambda_y I)^{-1}$$

Alternating Least Squares



$$U_x = r_x Y (Y^T Y + \lambda_x I)^{-1}$$

Alternating Least Squares

$$R = \begin{array}{ccccc} & \text{user 1} & \text{user 2} & \text{user 3} & \dots & \text{user N} & \\ \left[\begin{array}{c} 1 \\ 3.2 \\ 5 \\ \vdots \\ 2 \end{array} \right. & \begin{array}{c} 4.5 \\ 3 \\ 3 \\ \vdots \\ 4 \end{array} & \begin{array}{c} 3.8 \\ 3 \\ 3.4 \\ \vdots \\ 1 \end{array} & \begin{array}{c} \dots \\ \dots \\ \dots \\ \ddots \\ \dots \end{array} & \begin{array}{c} 3 \\ 4 \\ 3.1 \\ \vdots \\ 2.7 \end{array} & \begin{array}{c} \text{product 1} \\ \text{product 2} \\ \text{product 3} \\ \vdots \\ \text{product M} \end{array} \end{array}$$

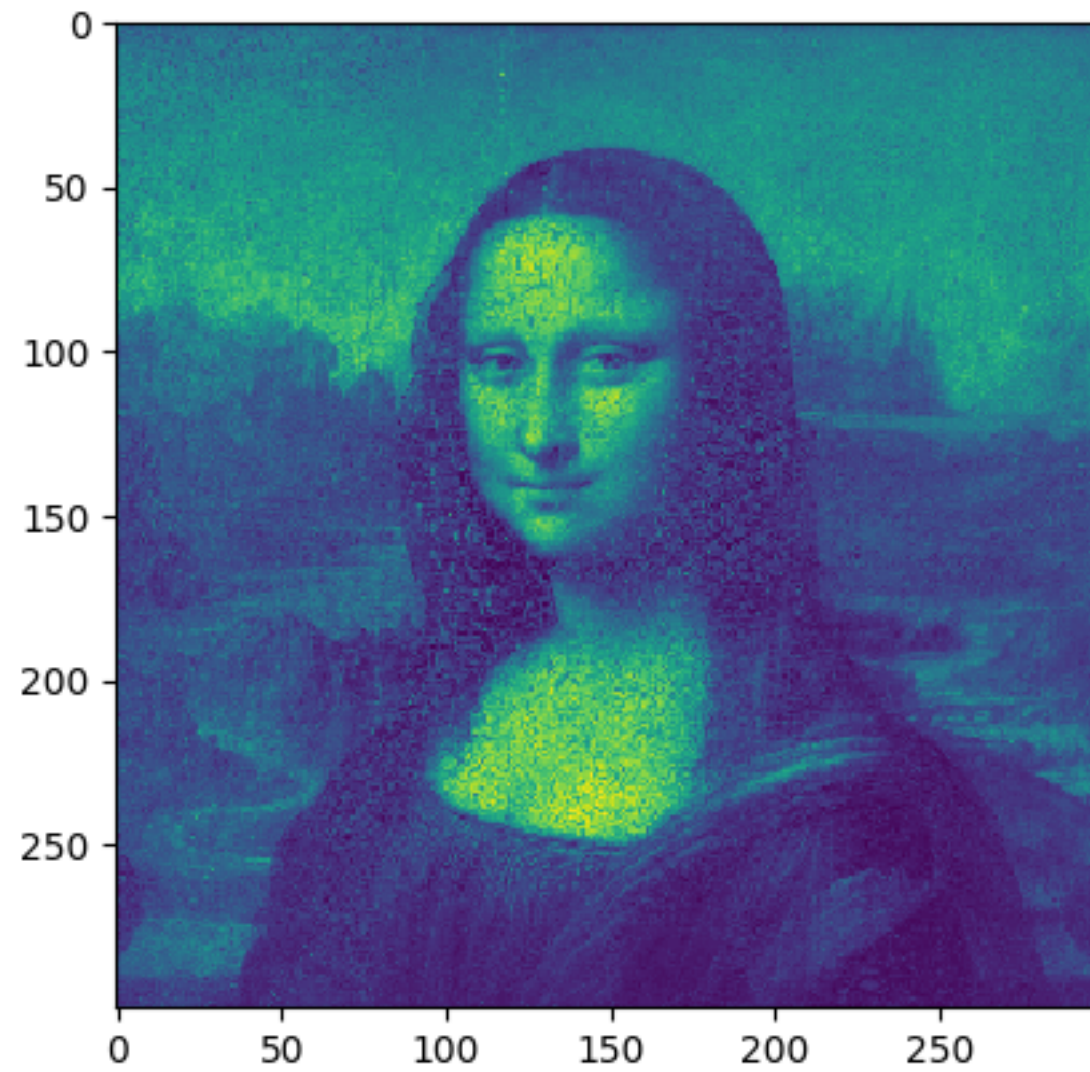
Batch ALS

	1	2	3	4	...	300
1	70	82	60	54		65
2	70	86	68	67		72
3	96	103	82	82		77
4	90	87	68	93		82
...						
300	38	48	44	51		35

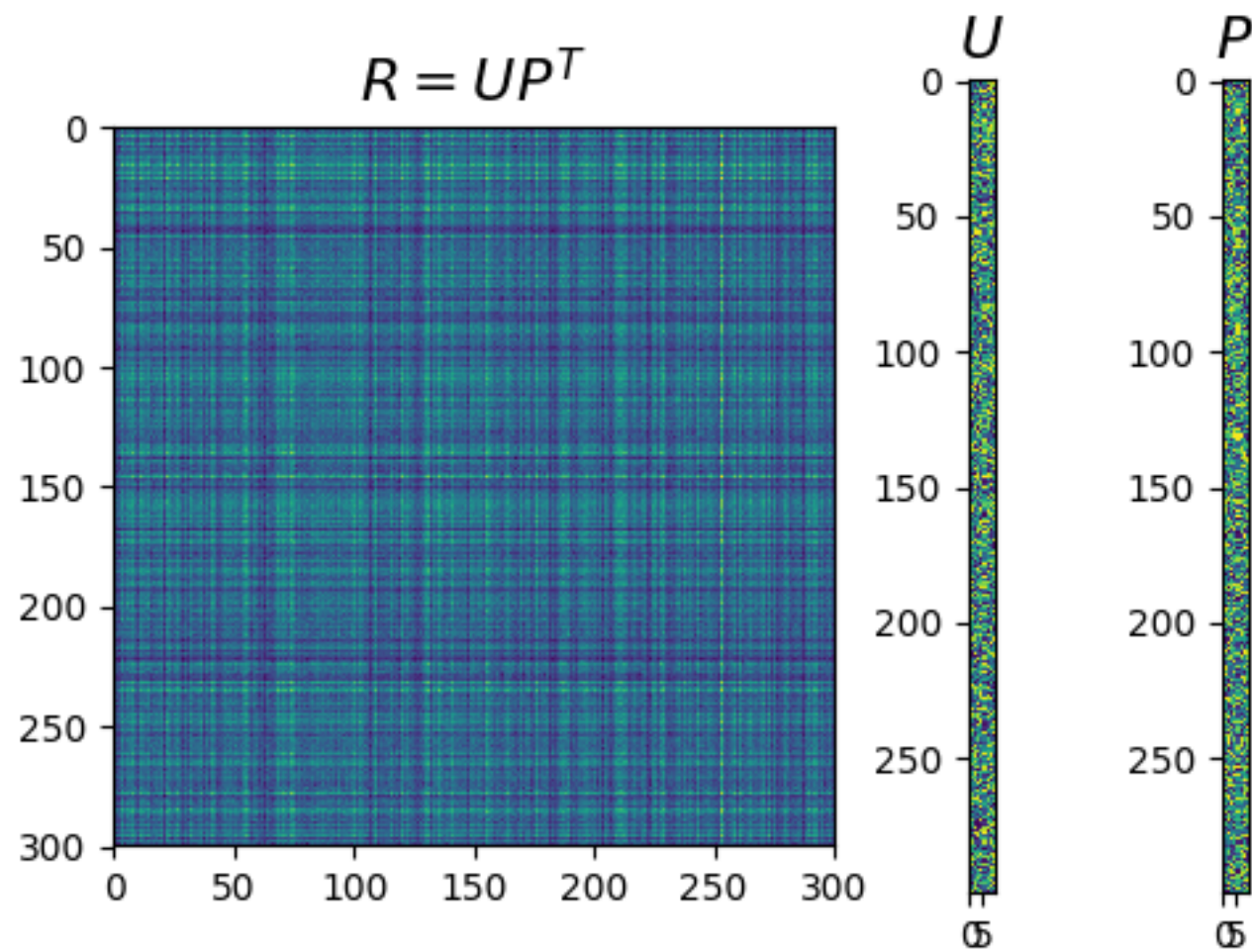
Batch ALS

	1	2	3	4	...	300
1	70	82	60	54		65
2	70	86	68	67		72
3	96	103	82	82		77
4	90	87	68	93		82
...						
300	38	48	44	51		35

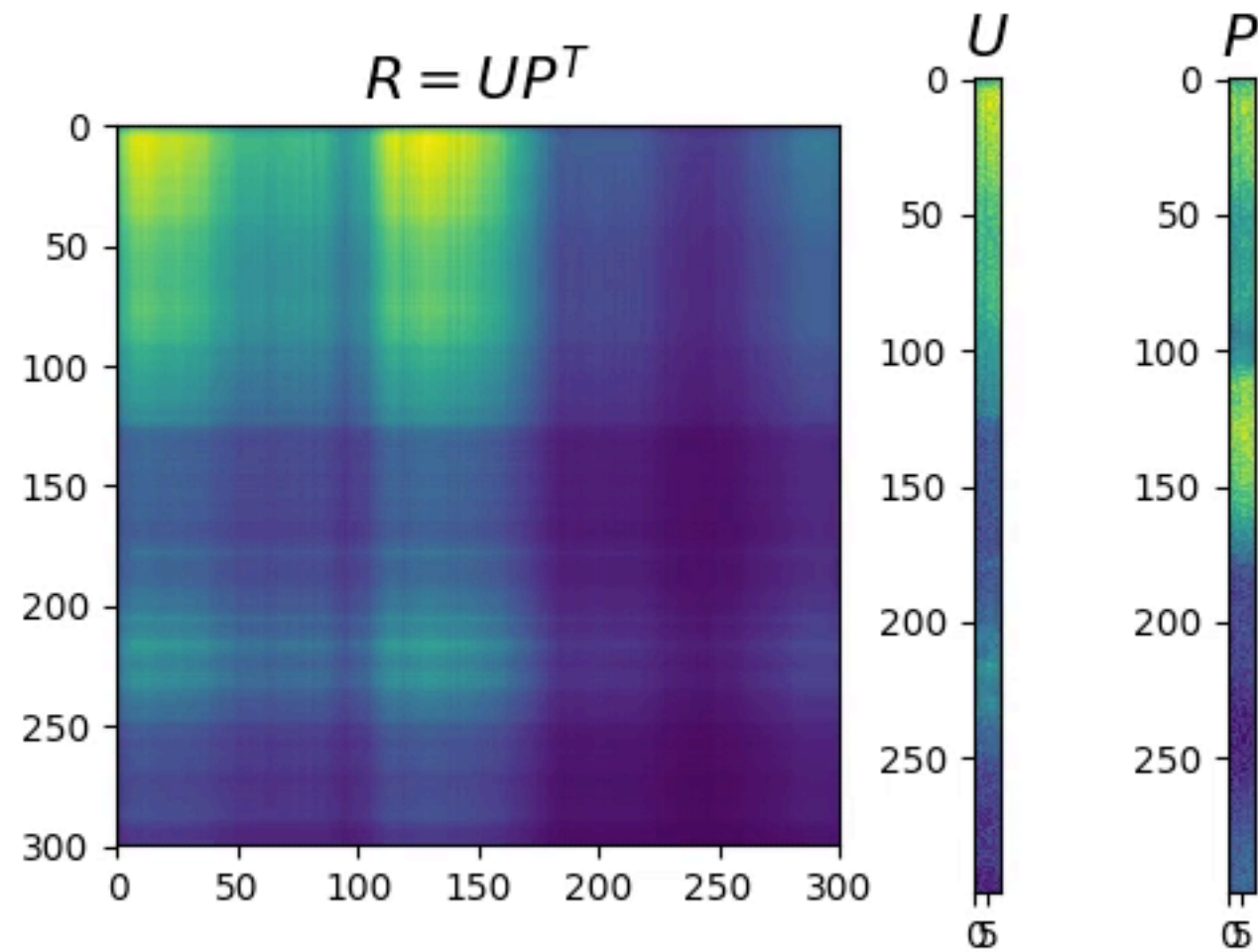
Batch ALS



Batch ALS



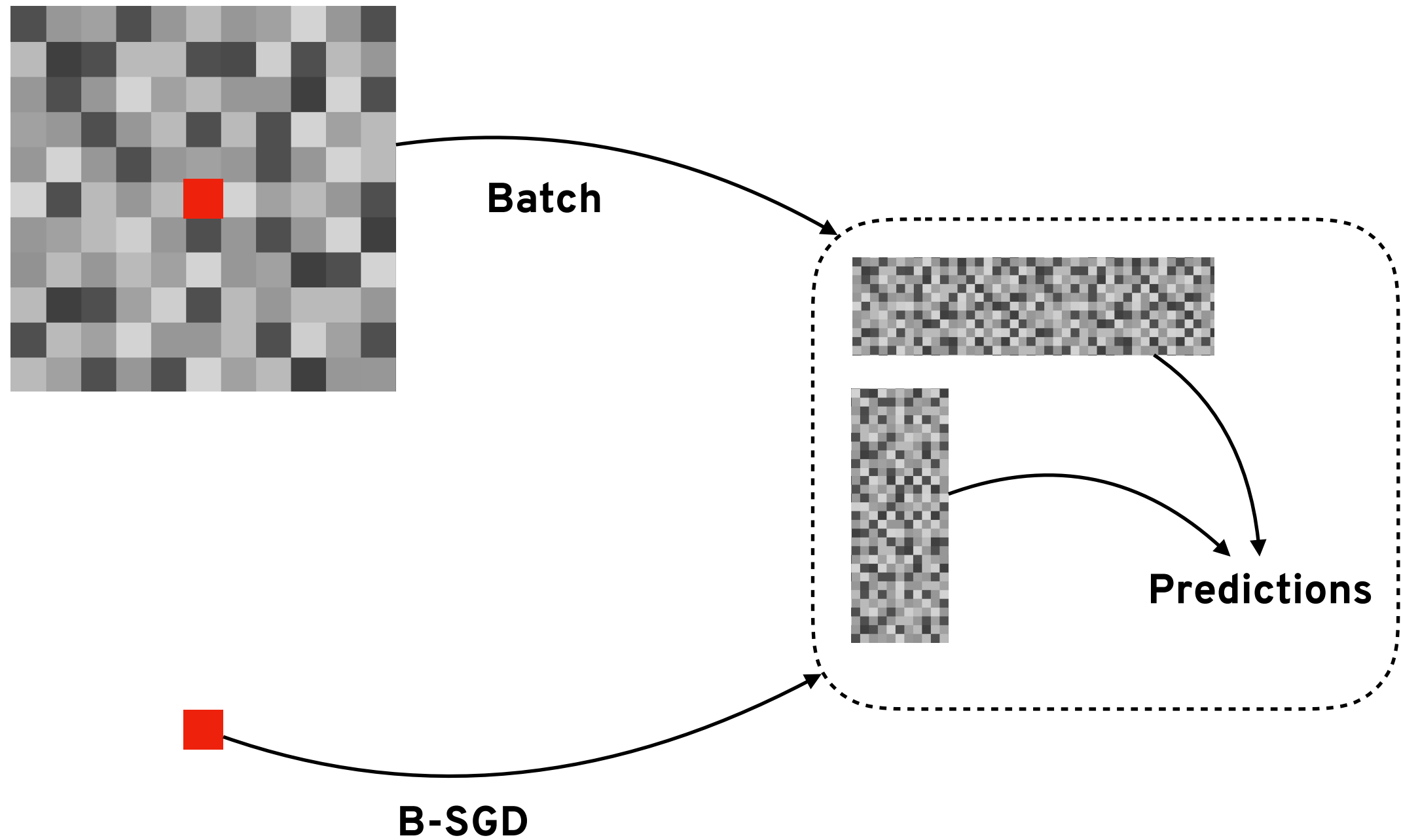
Batch ALS



Streaming ALS

- Can we update the model with a data stream?
- Stochastic Gradient Descent (SGD)
 - Bias SGD (B-SGD)

Streaming ALS



Streaming ALS

$$b_{x,y} = \mu + b_x + b_y$$

$$\hat{r}_{x,y} = \mu + b_x + b_y + \mathbf{U}_x \cdot \mathbf{P}_y^T$$

Streaming ALS

$$\hat{r}_{x,y} = \mu + b_x + b_y + \mathbf{U}_x \cdot \mathbf{P}_y^T$$

$$\text{loss} = \sum_{x,y} \left(\underbrace{r_{x,y} - \hat{r}_{x,y}}_{\epsilon_{x,y}} \right)^2 + \dots$$

Streaming ALS

bias

$$b_x \leftarrow b_x + \gamma (\epsilon_{x,y} - \lambda_x b_x)$$

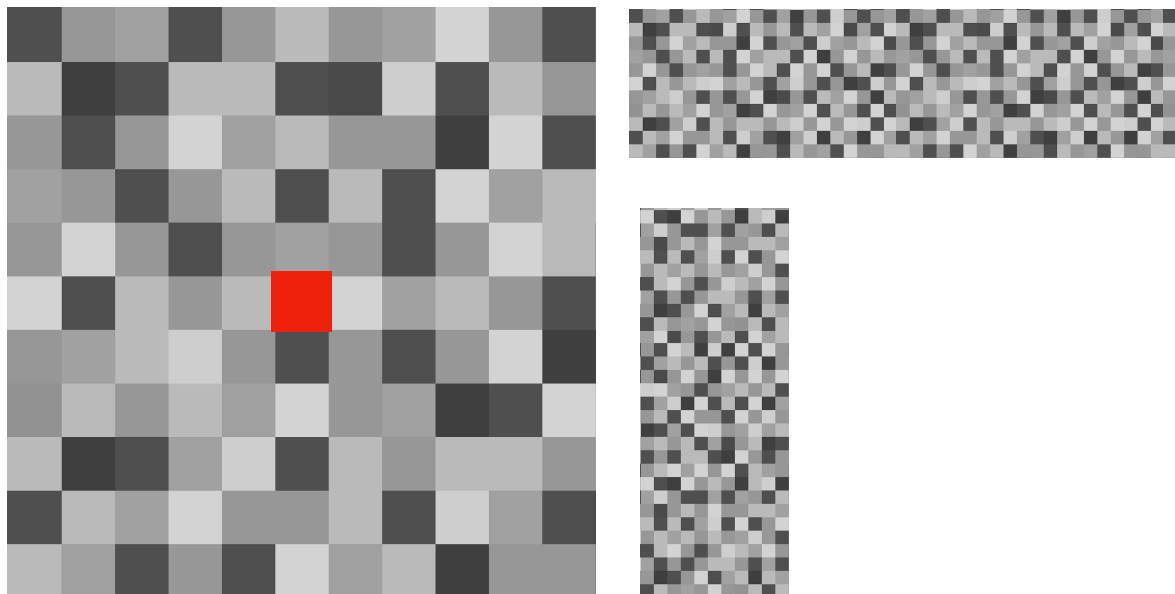
$$b_y \leftarrow b_y + \gamma (\epsilon_{x,y} - \lambda_y b_y)$$

factors

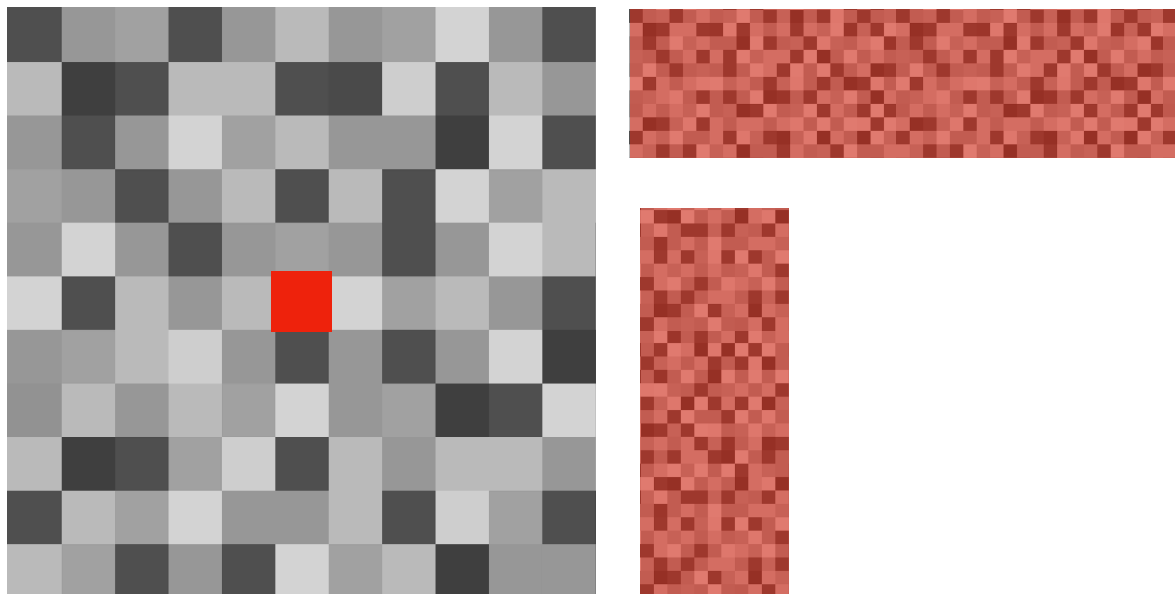
$$U_x \leftarrow U_x + \gamma (\epsilon_{x,y} P_y - \lambda'_x U_x)$$

$$P_y \leftarrow P_y + \gamma (\epsilon_{x,y} U_x - \lambda'_y P_y)$$

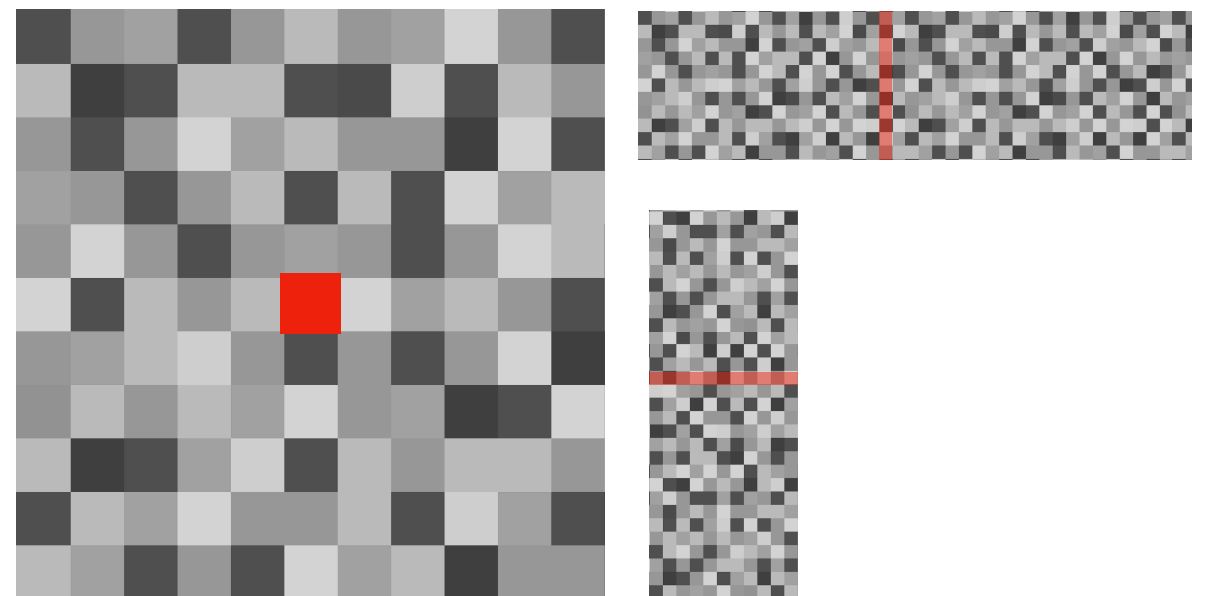
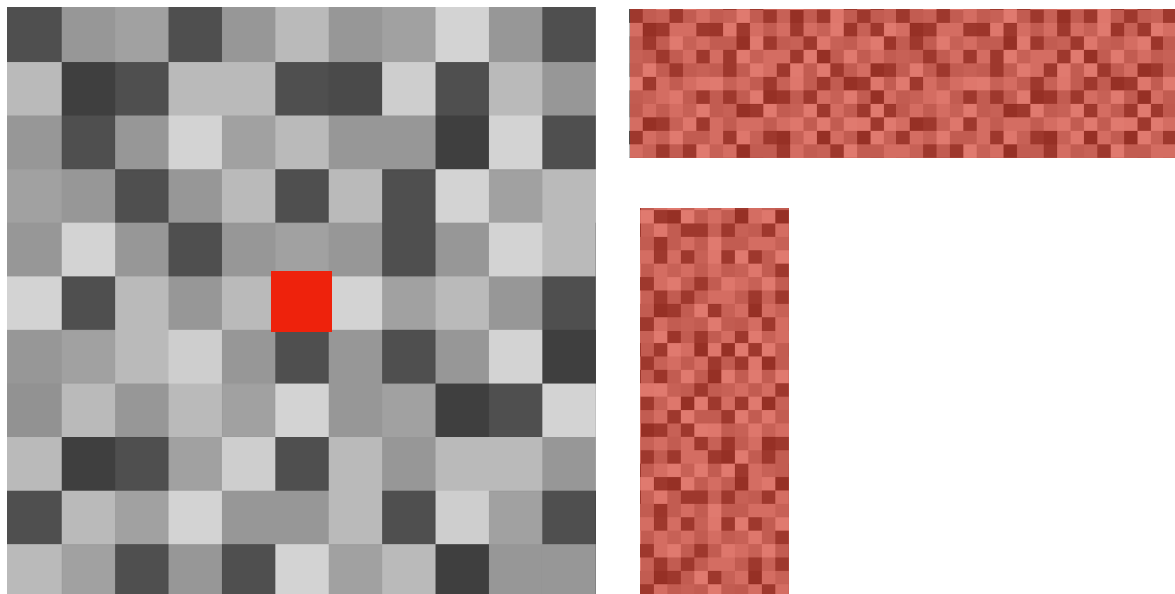
Streaming ALS



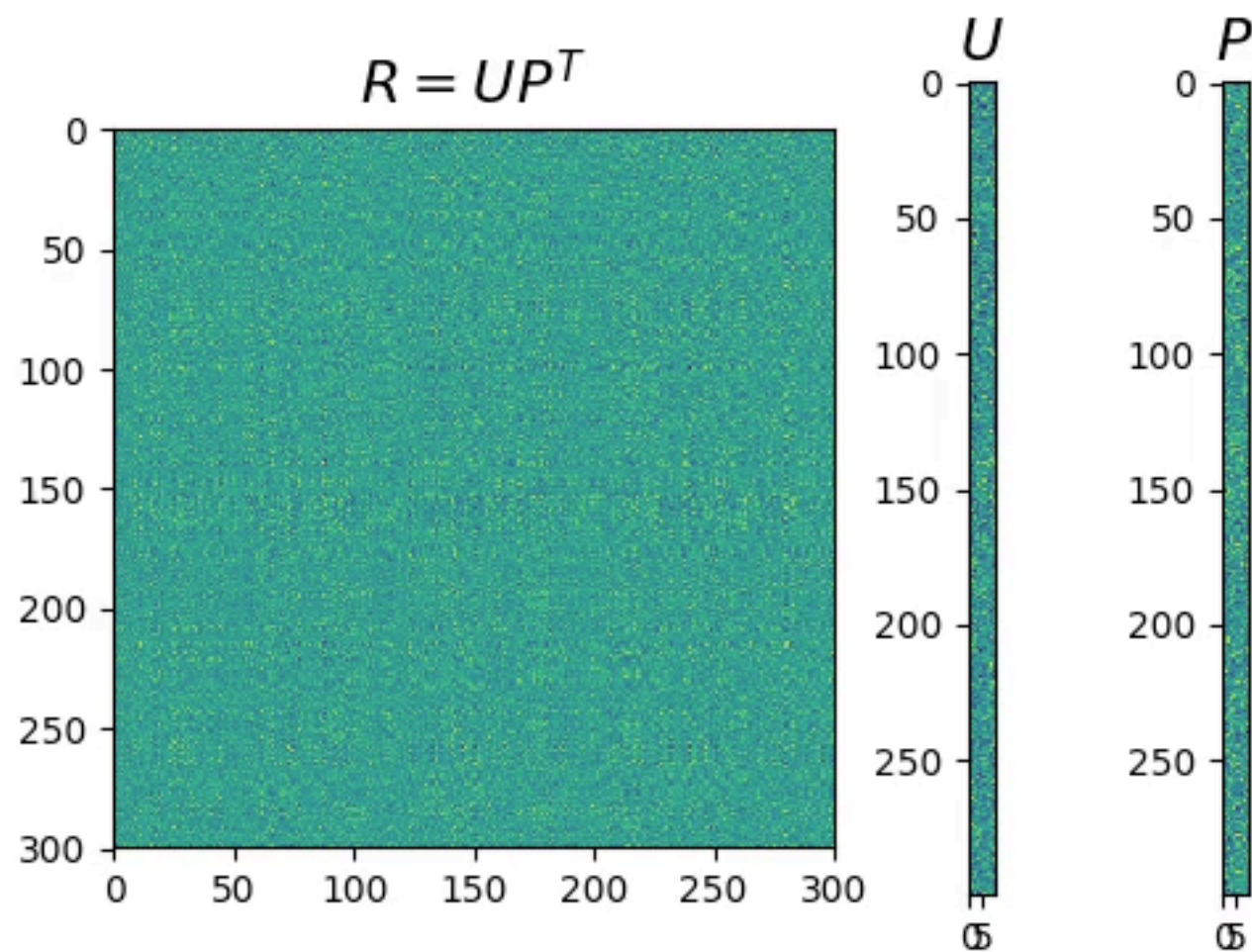
Streaming ALS



Streaming ALS



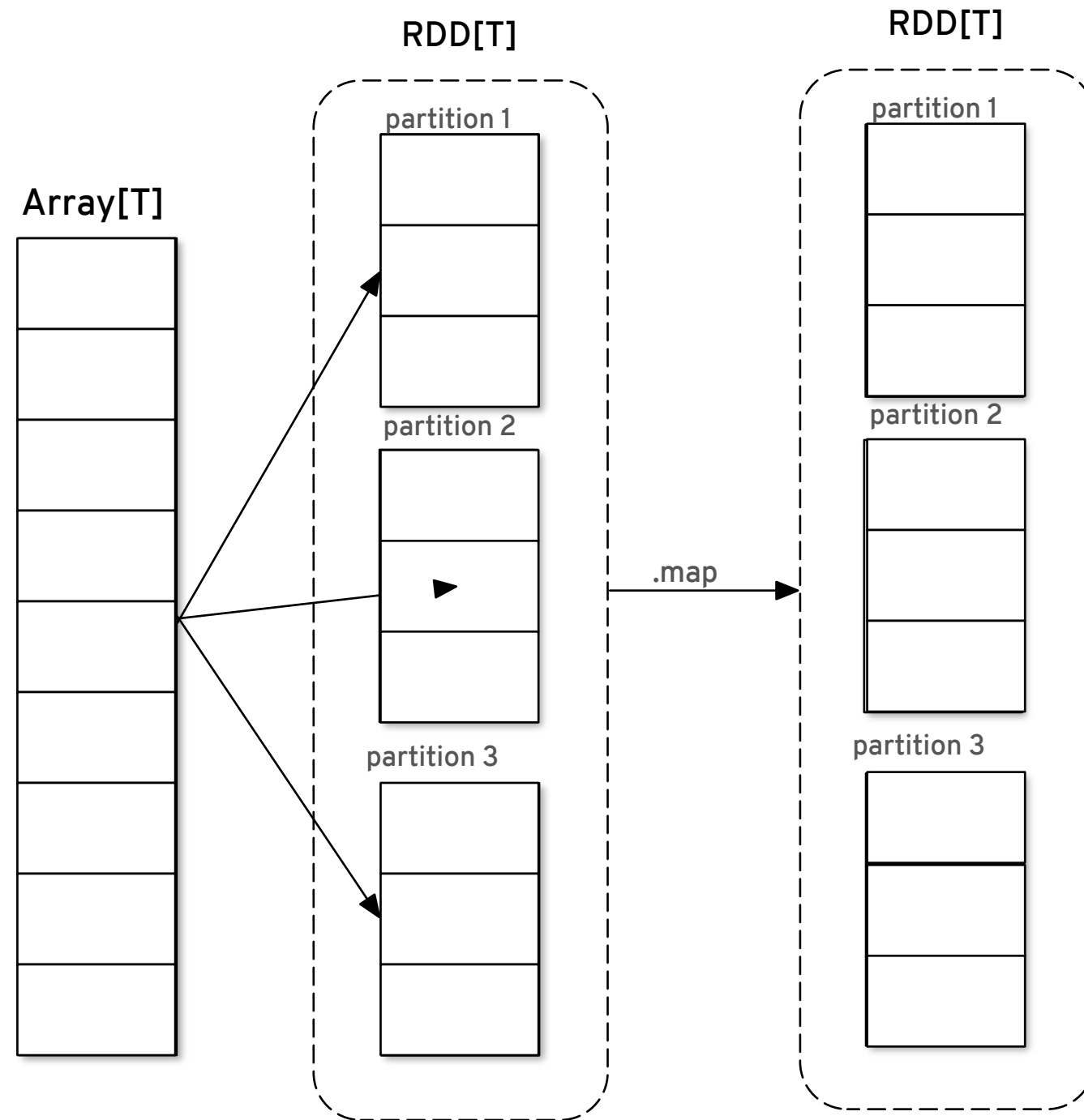
Streaming ALS



Apache Spark



Apache Spark



MLlib ALS

```
val model = ALS.train(ratings, rank, iterations, lambda)
```

MLlib ALS

```
val model = ALS.train(ratings, rank, iterations, lambda)
```

```
case class Rating(int user, int product, double rating)
```

```
val ratings: RDD[Rating]
```

MLlib ALS

```
val model = ALS.train(ratings, rank, iterations, lambda)
```

```
case class Rating(int user, int product, double rating)
```

```
val ratings: RDD[Rating]
```

```
val rank: int
```

MLlib ALS

```
val model = ALS.train(ratings, rank, iterations, lambda)
```

```
case class Rating(int user, int product, double rating)
```

```
val ratings: RDD[Rating]
```

```
val rank: int
```

```
val iterations: int
```


MLlib ALS

```
val model = ALS.train(ratings, rank, iterations, lambda)
```

```
case class Rating(int user, int product, double rating)
```

```
val ratings: RDD[Rating]
```

```
val rank: int
```

```
val iterations: int
```

```
val lambda: Double
```

MLlib ALS

```
> val model = ALS.train(ratings, rank, iterations, lambda)
```

```
model: MatrixFactorizationModel
```

```
class MatrixFactorizationModel {
```

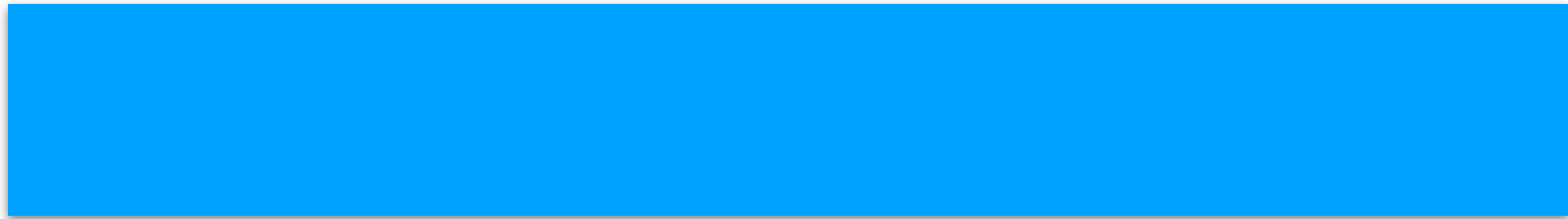
```
    val userFeatures: RDD[(Int, Array[Double])]
```

```
    val productFeatures: RDD[(Int, Array[Double])]
```

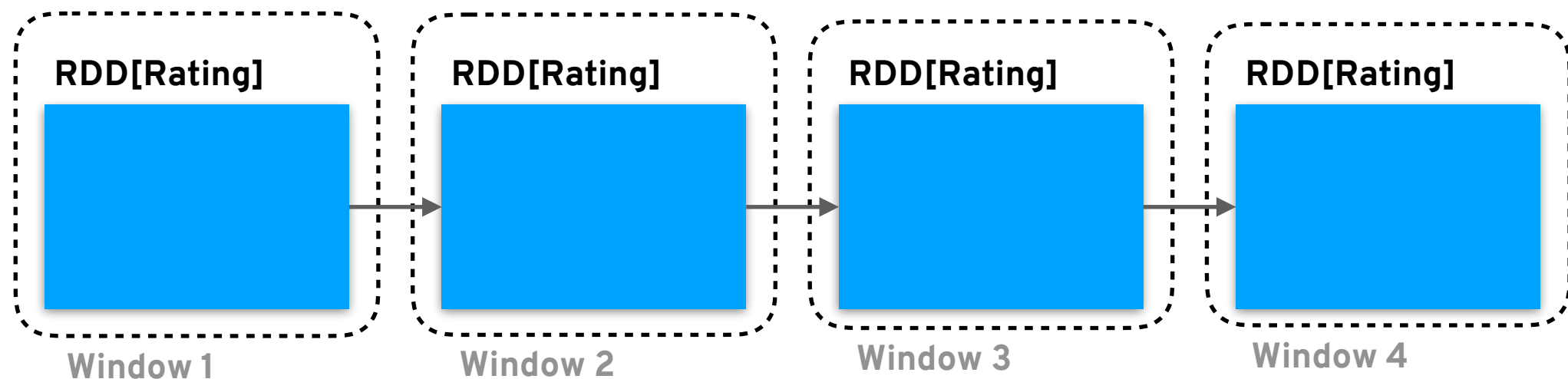
```
}
```

Spark Streaming ALS

RDD[Rating]



DStream[Rating]

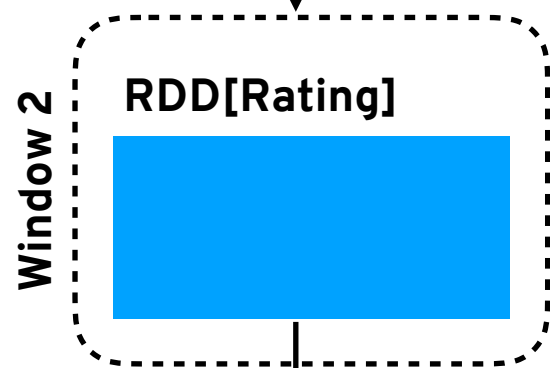


Spark Streaming ALS

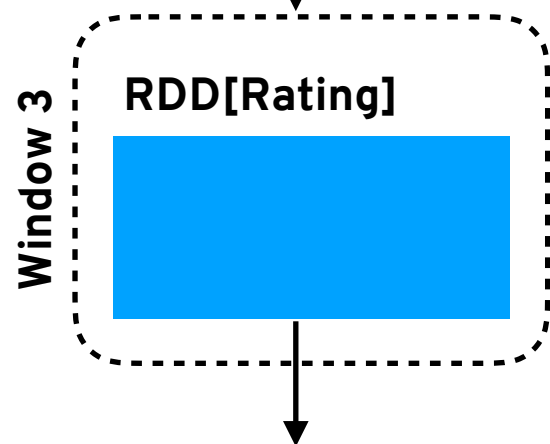
DStream[Rating]



```
model = StreamingALS.train(rdd1, params)
```



```
model = model.train(rdd2)
```



```
model = model.train(rdd3)
```

Spark Streaming ALS

```
userBias += gamma * (error - lambda * userBias)
```

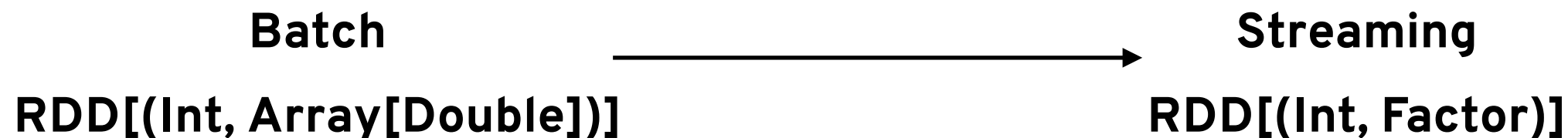
```
userFeature(i) += gamma * (error * prodFeature(j) - lambda * userFeature(i))
```

Spark Streaming ALS

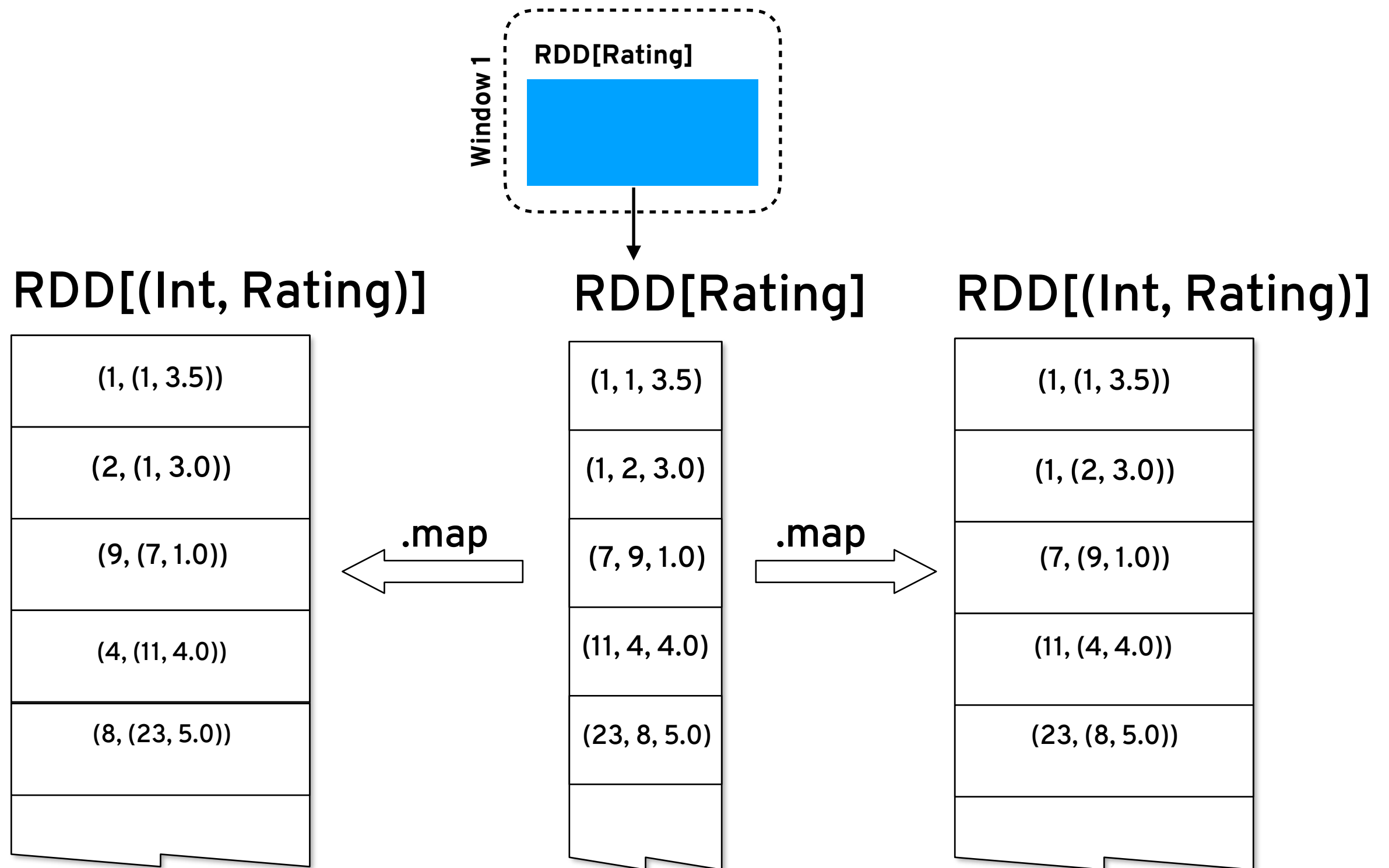
```
userBias += gamma * (error - lambda * userBias)
```

```
userFeature(i) += gamma * (error * prodFeature(j) - lambda * userFeature(i))
```

```
case class Factor(var bias: Double, features: Array[Double])  
  extends Serializable {  
  
}
```



Spark Streaming ALS



Spark Streaming ALS

RDD[(Int, Rating)]

(1, (1, 3.5))
(2, (1, 3.0))
(9, (7, 1.0))
(4, (11, 4.0))
(8, (23, 5.0))

.map
→

RDD[(Int, Factor)]

(1, [0.123, -0.234, ...])
(2, [0.943, 0.527, ...])
(9, [0.421, -0.594, ...])
(4, [0.034, 0.661, ...])
(8, [0.713, -0.335, ...])

Spark Streaming ALS

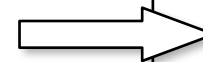
.join

RDD[(Int, Rating)]

(1, (1, 3.5))
(2, (1, 3.0))
(9, (7, 1.0))
(4, (11, 4.0))
(8, (23, 5.0))

RDD[(Int, Factor)]

(1, [0.123, -0.234, ...])
(2, [0.943, 0.527, ...])
(9, [0.421, -0.594, ...])
(4, [0.034, 0.661, ...])
(8, [0.713, -0.335, ...])



RDD[(Int, (Int, Double, Factor))]

(1, (1, 3.5, [0.123, ...]))
(2, (1, 3.0, [0.943, ...]))
(9, (7, 1.0, [0.421, ...]))
(4, (11, 4.0, [0.034, ...]))
(8, (23, 5.0, [0.713, ...]))

Spark Streaming ALS

RDD[(Int, (Int, Double, Factor), Factor)]

(1, (1, 3.5, [0.123, ...]), [0.426, ...])

$$\hat{r}_{x,y} = \mu + b_x + b_y + U_x \cdot P_y^T$$

```
prediction = dot(userFactors.features, itemFactors.features)
+ userFactors.bias + itemFactors.bias + bias
```

Spark Streaming ALS

RDD[(Int, (Int, Double, Factor), Factor)]

(1, (1, 3.5, [0.123, ...]), [0.426, ...])

$$\hat{r}_{x,y} = \mu + b_x + b_y + U_x \cdot P_y^T$$

```
prediction = dot(userFactors.features, itemFactors.features)
+ userFactors.bias + itemFactors.bias + bias
```

$$\epsilon_{x,y} = r_{x,y} - \hat{r}_{x,y}$$

```
eps = rating - prediction
```

Spark Streaming ALS

RDD[(Int, (Int, Double, Factor), Factor)]

(1, (1, 3.5, [0.123, ...]), [0.426, ...])

$$\hat{r}_{x,y} = \mu + b_x + b_y + U_x \cdot P_y^T$$

```
prediction = dot(userFactors.features, itemFactors.features)
+ userFactors.bias + itemFactors.bias + bias
```

$$\epsilon_{x,y} = r_{x,y} - \hat{r}_{x,y}$$

```
eps = rating - prediction
```

$$\gamma (\epsilon_{x,y} U_x - \lambda'_y P_y)$$

```
(0 until rank).map { i =>
  gamma * (eps * userFactors.features(i) - lambda * itemFactors.features(i))
}
```

Spark Streaming ALS

RDD[(Int, (Int, Double, Factor), Factor)]

(1, (1, 3.5, [0.123, ...]), [0.426, ...])

$$\hat{r}_{x,y} = \mu + b_x + b_y + U_x \cdot P_y^T$$

```
prediction = dot(userFactors.features, itemFactors.features)
+ userFactors.bias + itemFactors.bias + bias
```

$$\epsilon_{x,y} = r_{x,y} - \hat{r}_{x,y}$$

```
eps = rating - prediction
```

$$\gamma (\epsilon_{x,y} U_x - \lambda'_y P_y)$$

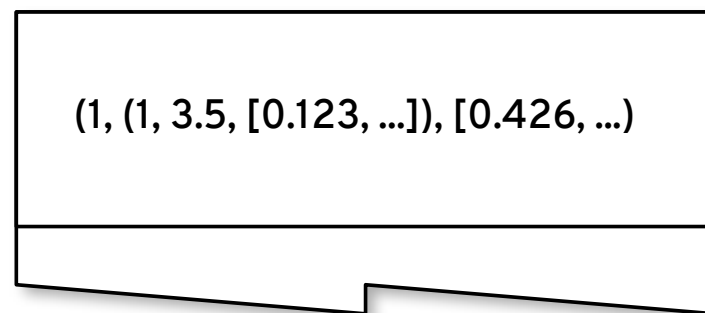
```
(0 until rank).map { i =>
  gamma * (eps * userFactors.features(i) - lambda * itemFactors.features(i))
}
```

$$\gamma (\epsilon_{x,y} - \lambda_y b_y)$$

```
gamma * (eps - lambda * itemFactors.bias)
```

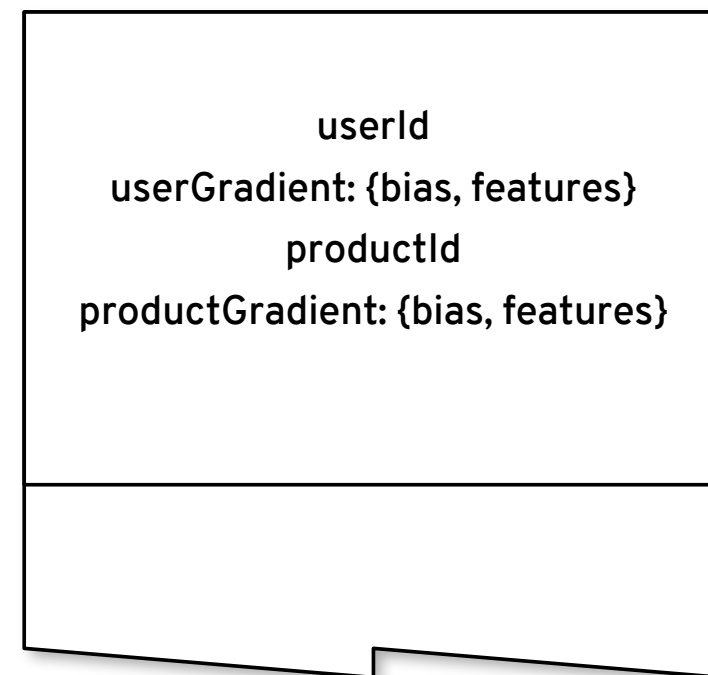
Spark Streaming ALS

RDD[(Int, (Int, Double, Factor), Factor)]



`.map`
→

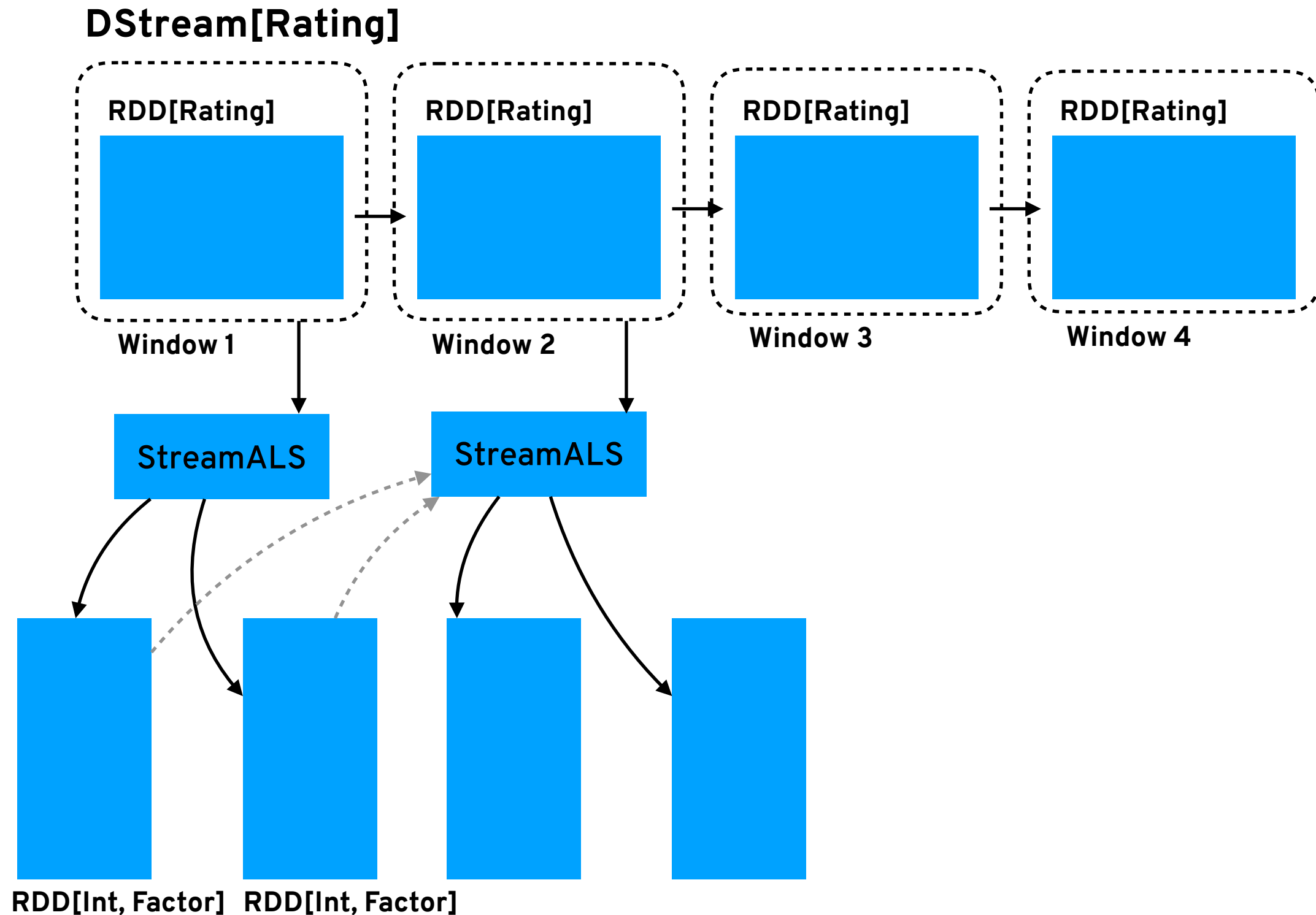
RDD[(Int, Factor, Int, Factor)]



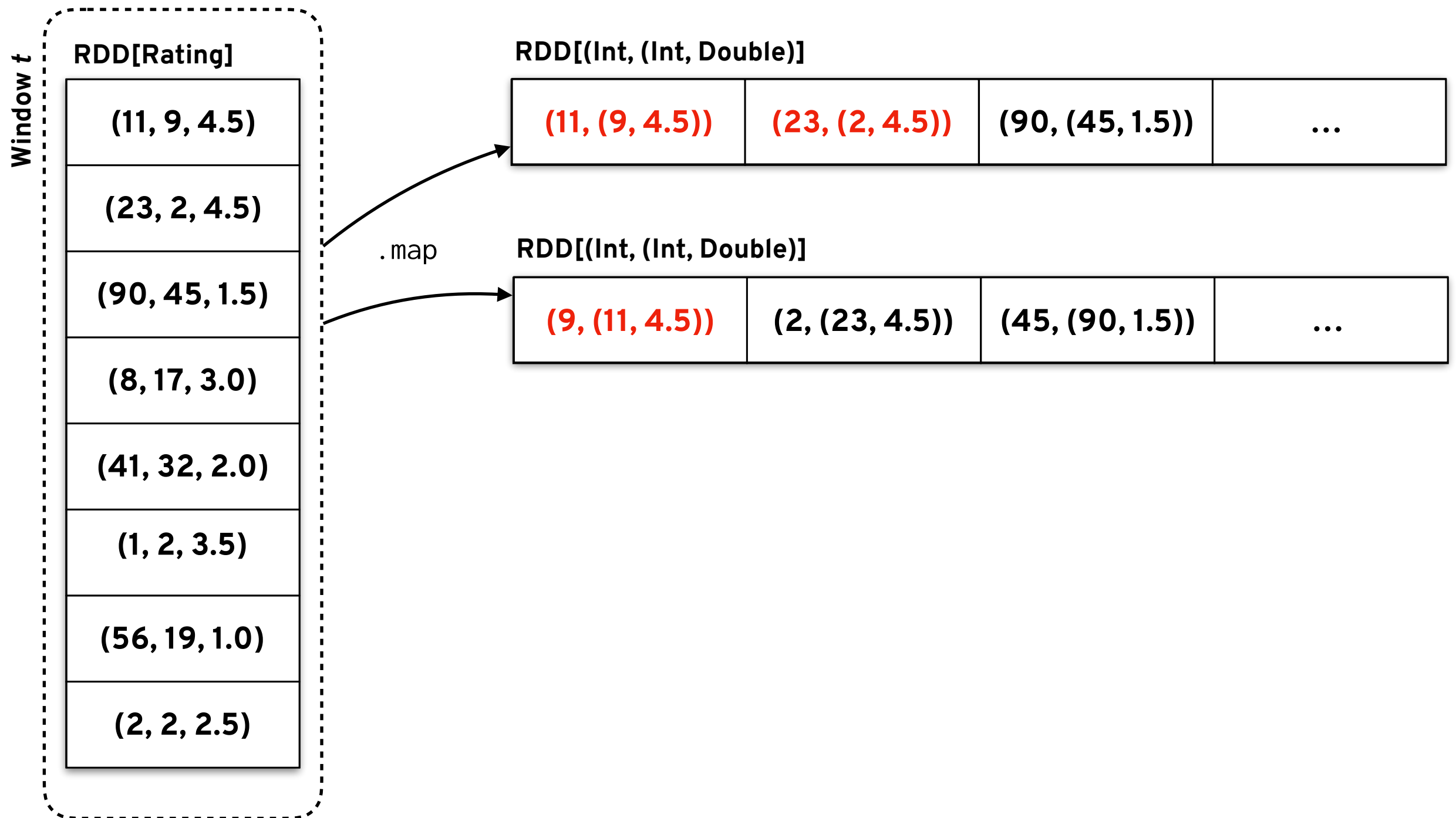
Spark Streaming ALS

- We need to `.aggregateByKey`
 - Bias, features for each user and product
- Finally recalculate final bias, features using `.leftJoin`
 - Add “original” Factor to gradient

Spark Streaming ALS



Spark Streaming ALS



Spark Streaming ALS

RDD[(Int, (Int, Double))]

(11, (9, 4.5))	(23, (2, 4.5))	(90, (45, 1.5))	...
----------------	----------------	-----------------	-----

RDD[(Int, (Int, Factor))]

(11, (0.12, [0.9,...]))
(23, (-0.1, [1.37,...]))
(1, (1, [0.123,...]))
...

.fullOuterJoin

(11, (9, 4.5), 11, (0.12, [0.9,...]))
(23, (2, 4.5), (23, (-0.1, [1.37,...])))
(90, (45, 1.5), None)
...

```
userRatings.fullOuterJoin(userFactors).map {  
  case (userId, (_, userFactors)) =>  
    (userId, userFactors(featureGenerator.nextValue()))  
}
```

RDD[(Int, (Int, Double), Int, Factor)]

Data

- MovieLens
- Widely used in recommendation engine research
- Variants
 - Small - 100,000 ratings / 9,000 movies / 700 users
 - Full - 26 million ratings / 45,000 movies / 270,000 users
- CSV data
 - Ratings
 - (userId, movieId, rating, timestamp)
 - (100, 200, 3.5, 2010-12-10 12:00:00)

Training batch ALS

```
val split: Array[RDD[Rating]] = ratings.randomSplit(0.8, 0.2)  
val model = ALS.train(split(0), rank, iter, lambda)
```

Training batch ALS

```
val split: Array[RDD[Rating]] = ratings.randomSplit(0.8, 0.2)
```

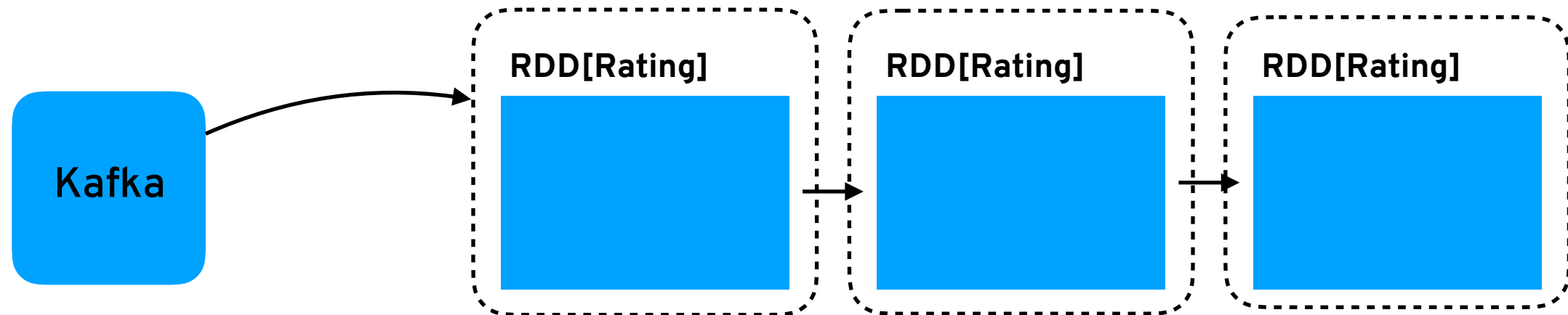
```
val model = ALS.train(split(0), rank, iter, lambda)
```

```
val predictions: RDD[Rating] = model.predict(split(1).map { x =>  
  (x.user, x.product)  
})
```

```
val pairs = predictions.map(x => ((x.user, x.product), x.rating))  
  .join(split(1).map(x => ((x.user, x.product), x.rating))  
  .values
```

```
Val RMSE = math.sqrt(pairs.map(x => math.pow(x._1 - x._2, 2)).mean())
```

Training streaming ALS



```
val model = StreamingALS(rank, iterations, lambda, gamma)

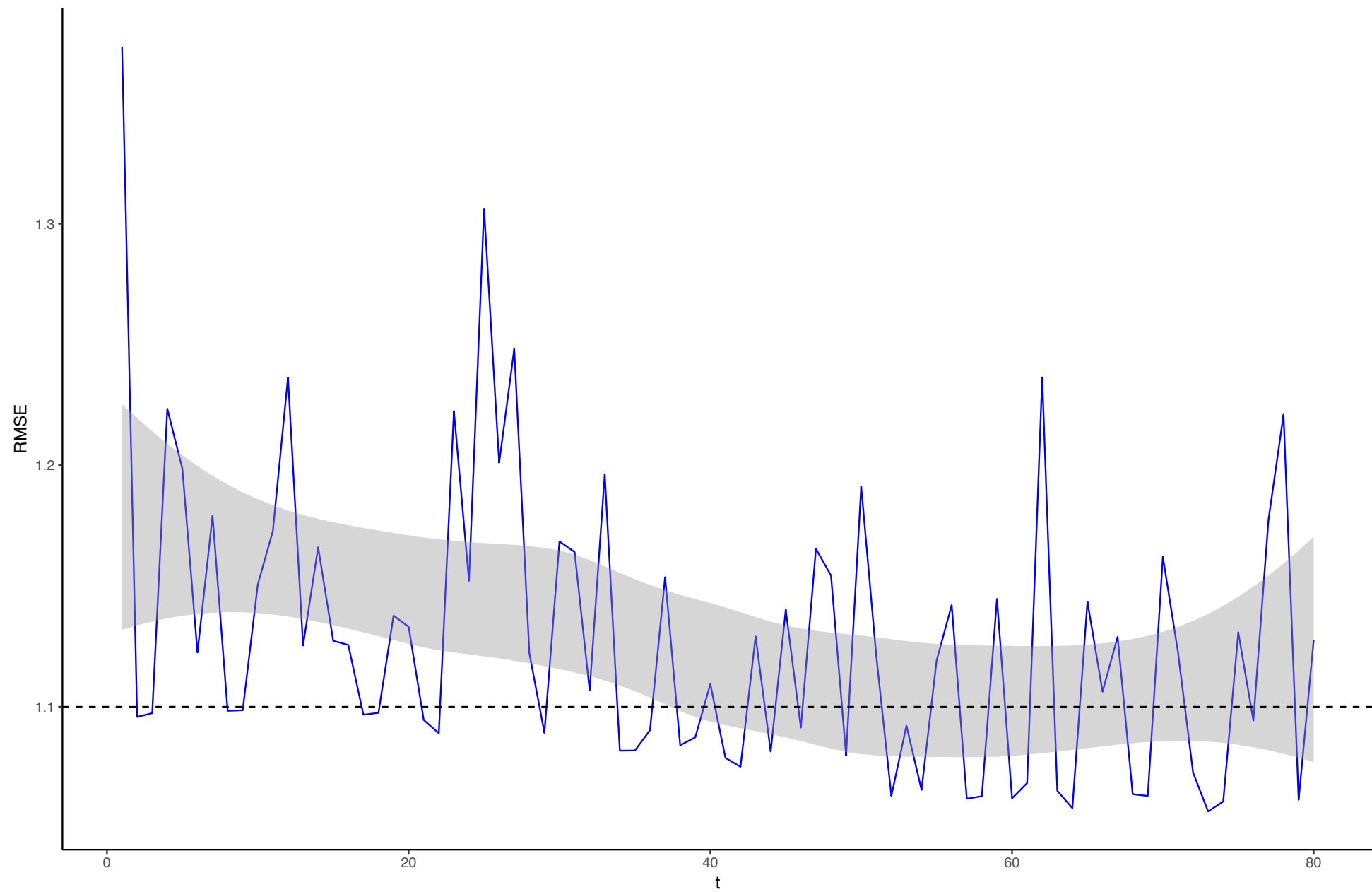
trainingStreamSet.foreachRDD { rdd =>

    model.train(rdd)

    val RMSE = calculateRMSE(model, validation)

}
```

Comparison

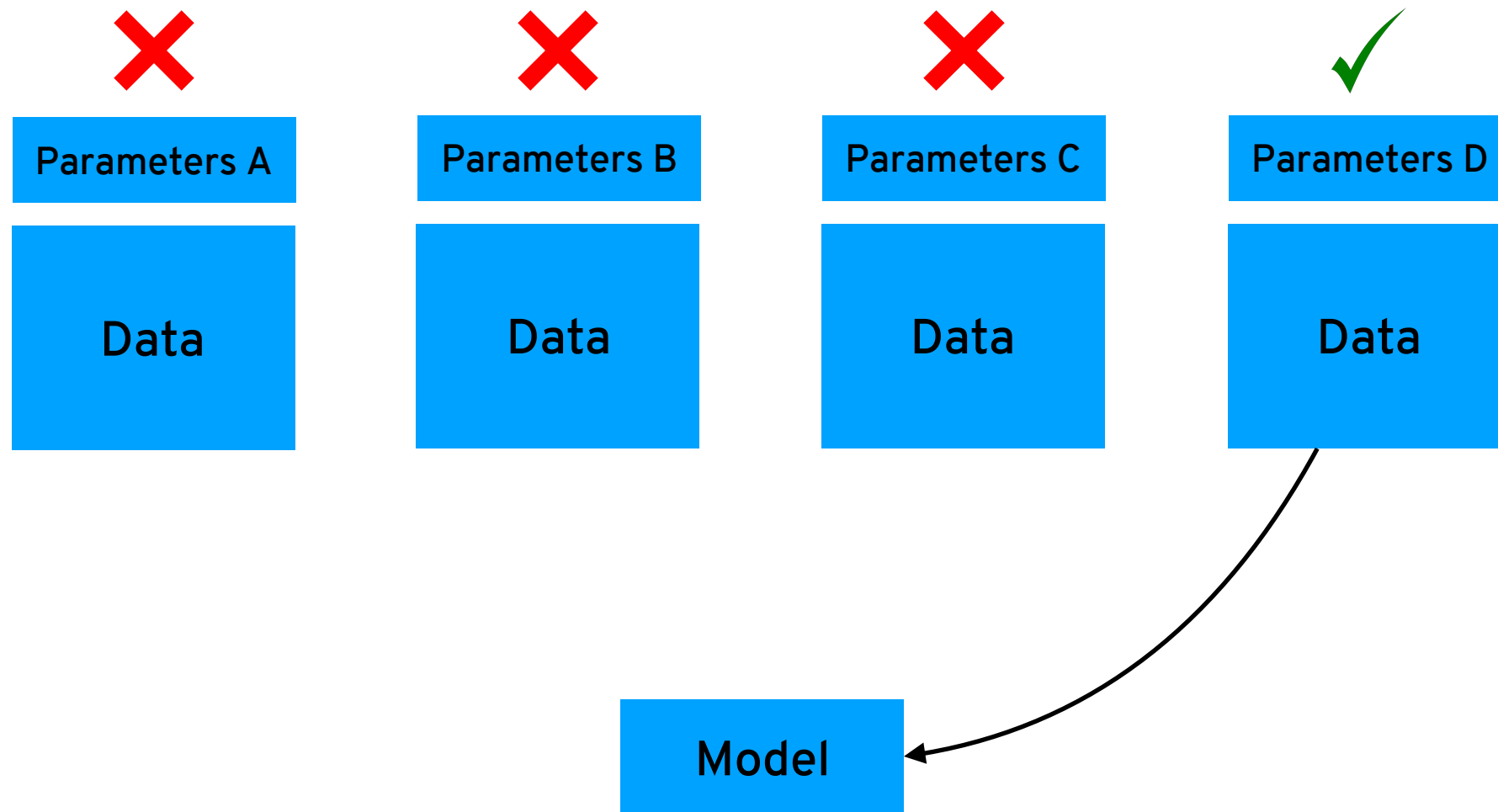


To consider

- **“Cold start”**
- **Same as batch ALS**
- **Too few observations = meaningless**
- **Train offline**

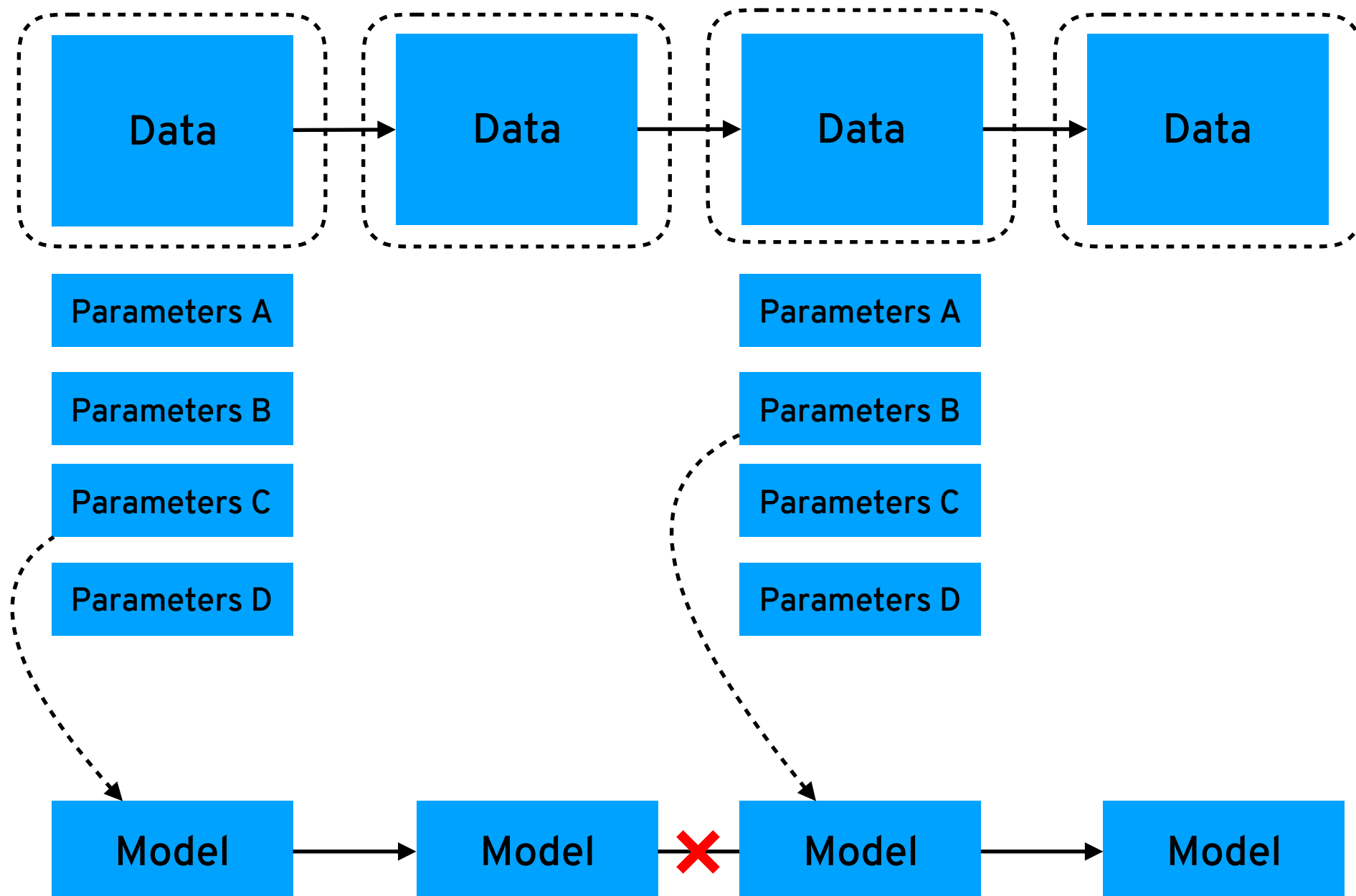
To consider

- Hyper-parameter estimation



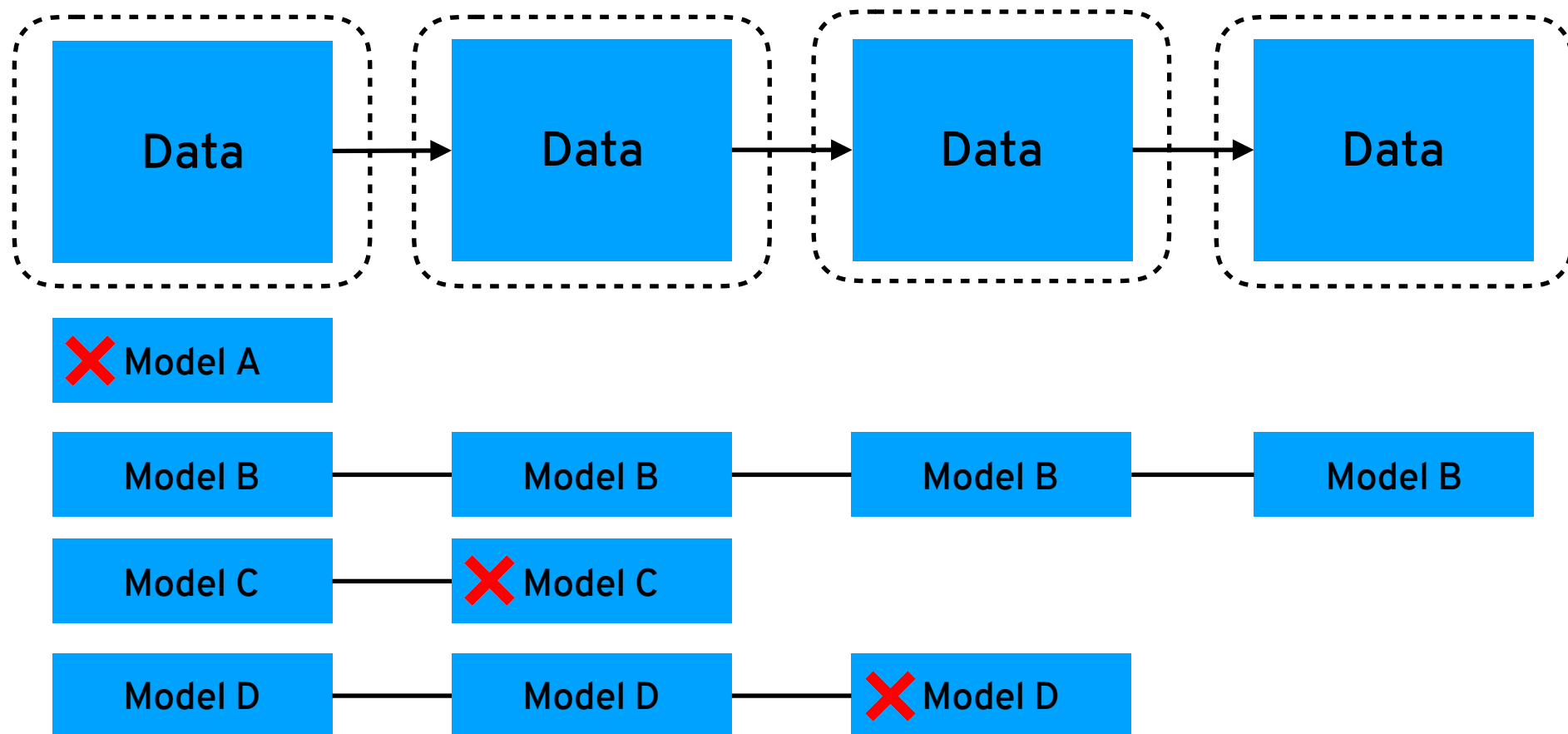
To consider

- Hyper-parameter estimation



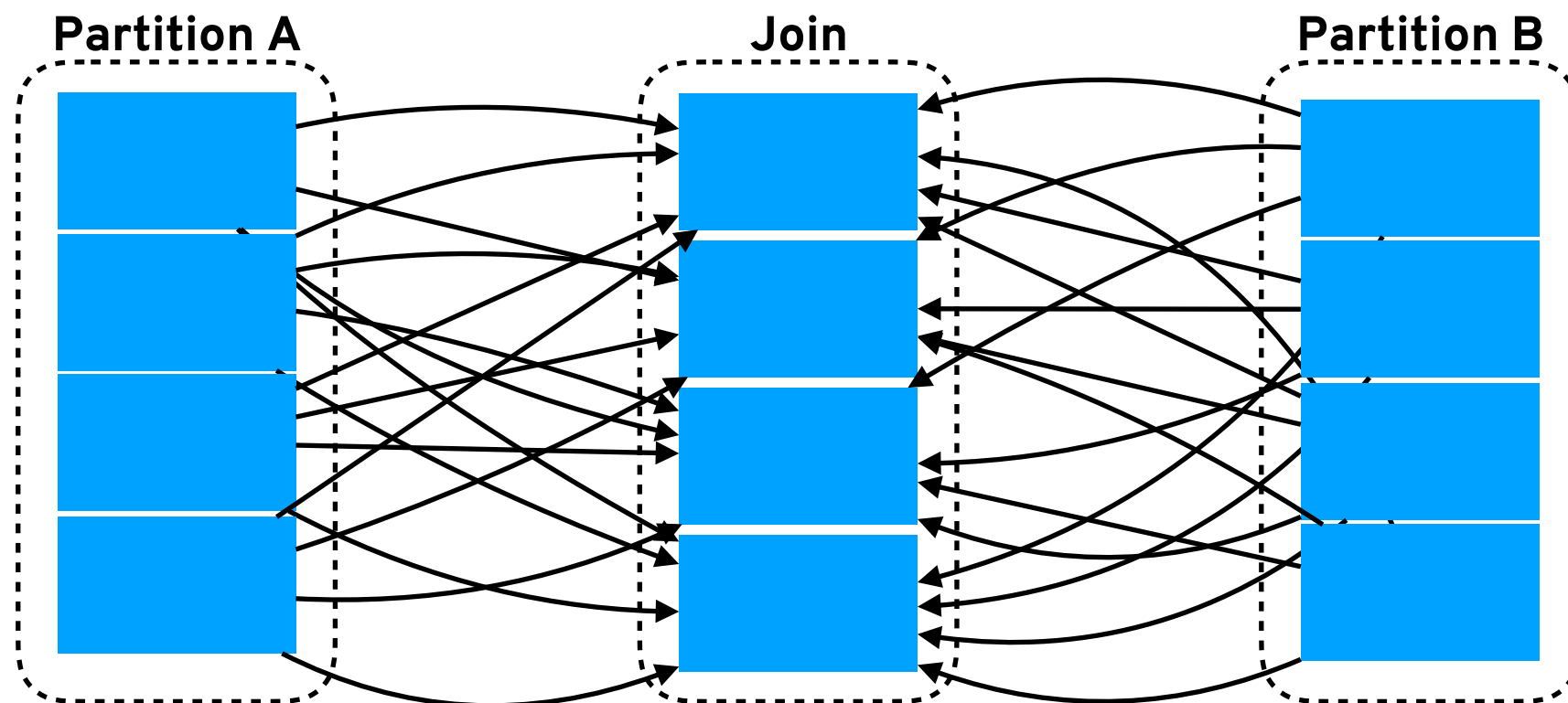
To consider

- Hyper-parameter estimation



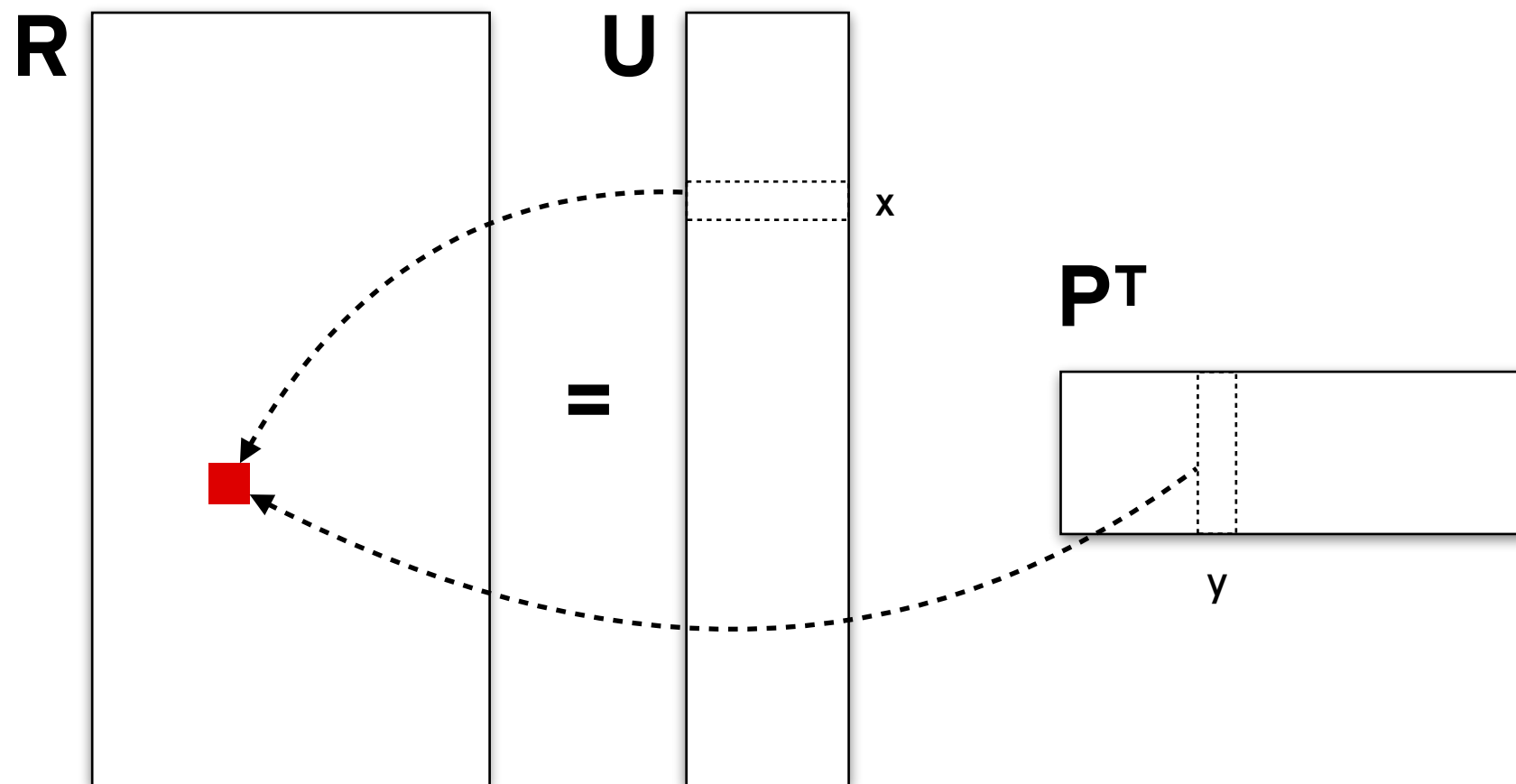
To consider

Partitioning



To consider

RDD random access?



```
val u = userFeatures.lookup(userId)
val p = productFeatures.lookup(productId)
val predicted = model.predict(userId, productId, u, p, globalBias)
```

Thank you!