# Dimension-Free Parameterized Approximation Schemes for Hybrid Clustering

## Ameet Gadekar

CISPA Helmholtz Center for Information Security
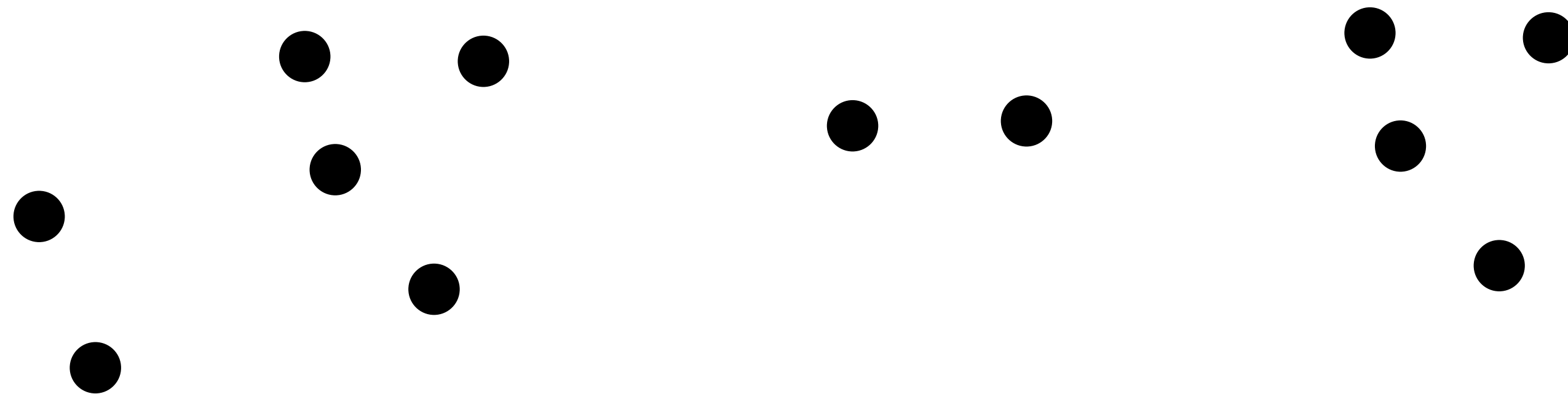Saarbrücken, Germany

## Tanmay Inamdar

Indian Institute of Technology Jodhpur
Jodhpur, India

07.03.2025     STACS 2025     Jena, Germany

# Clustering

Given a set of objects

Want to group them such that

objects in the same group are more "similar" to each other than to those in the other groups

# Clustering

Given a set of objects

Want to group them such that

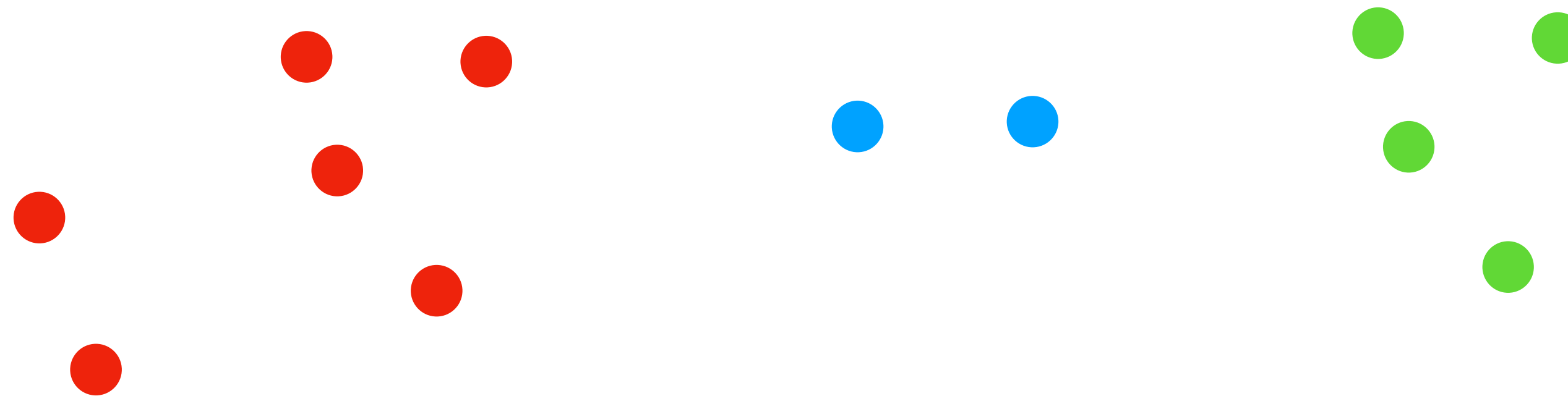objects in the same group are more "similar" to each other than to those in the other groups

# Clustering

Given a set of objects

Want to group them such that

objects in the same group are more "similar" to each other than to those in the other groups

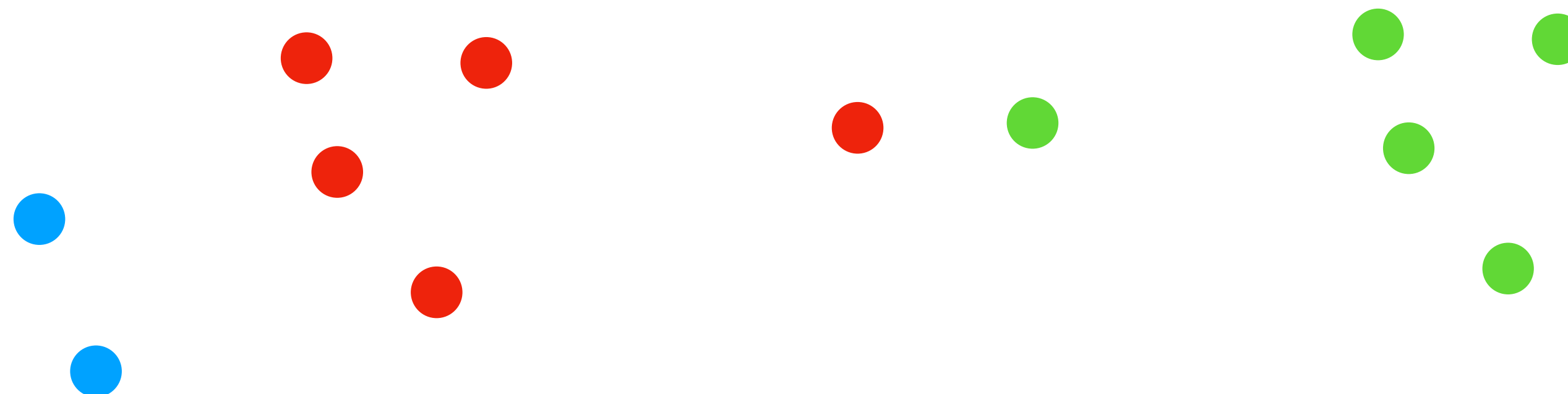Groups are called Clusters

# Center-based Clustering

Set of Objects — Points/Clients

Set of potential centers — Facilities

Want to choose centers and assign every point to a closest center

to minimize a clustering objective

# Center-based Clustering

Set of Objects — Points/Clients

Set of potential centers — Facilities

Want to choose centers and assign every point to a closest center

to minimize a clustering objective

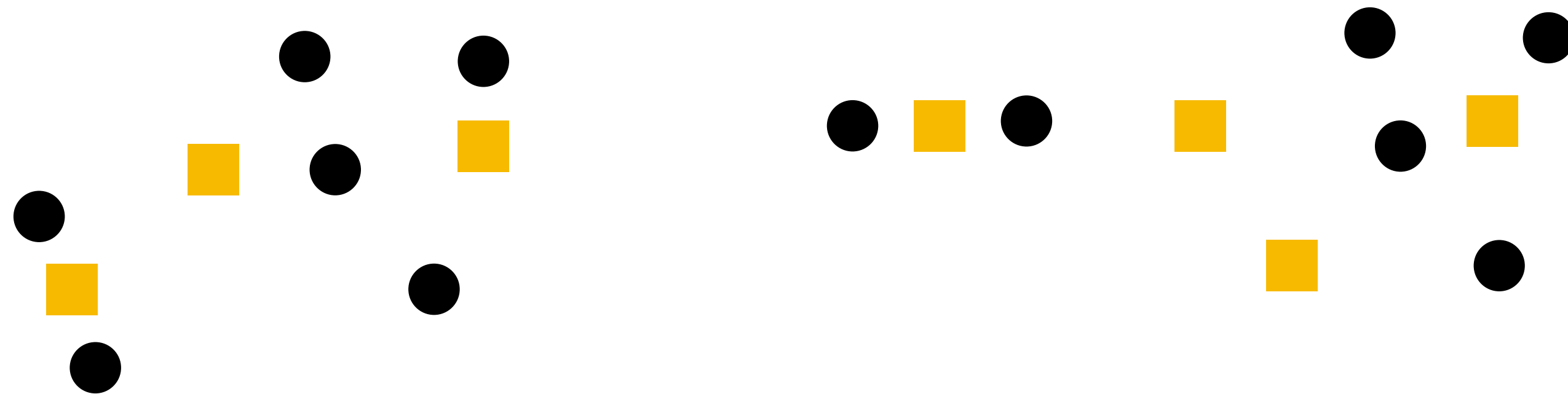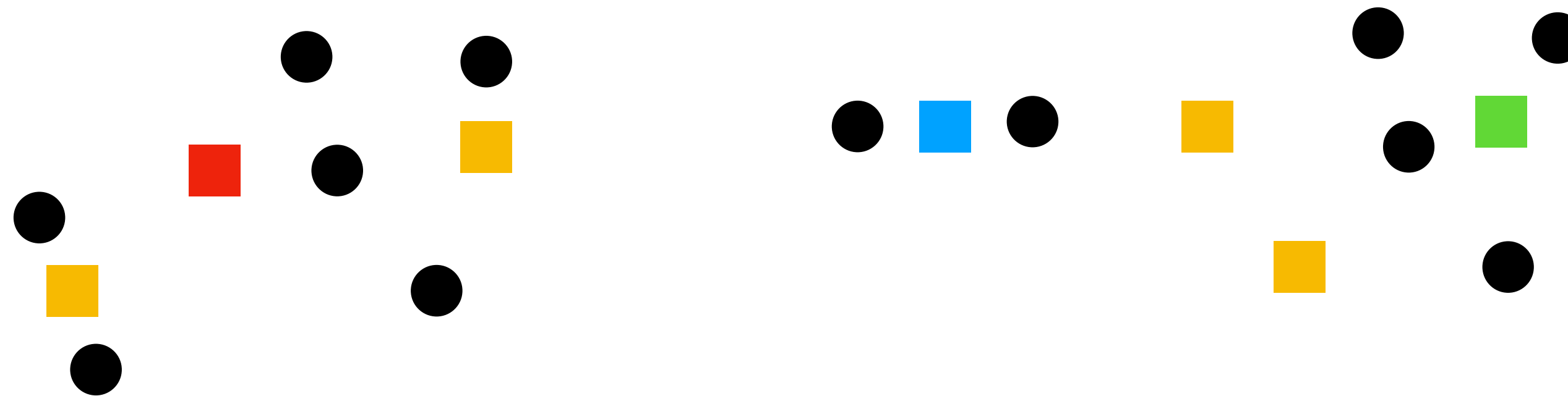# Center-based Clustering

Set of Objects — Points/Clients

Set of potential centers — Facilities

Want to choose centers and assign every point to a closest center

to minimize a clustering objective

# Center-based Clustering

- Input

    $P$: set of $n$ points

    $F$: set of facilities

    $d$: distance function on $P \cup F$

    $k$: positive integer

- Output

    $X \subseteq F$: set of $k$ centers

- Minimize an objective

    $k$-Median:  $\sum_{p \in P} d(p, X)$

    $k$-Center:  $\max_{p \in P} d(p, X)$

    $k$-Means:  $\sum_{p \in P} d(p, X)^2$

# Hybrid Clustering

- Interpolates between $k$-Median and $k$-Center

- Think of placing $k$ WiFi routers, each with coverage radius $r$

- Clients within coverage, pay 0 (zero)

- Clients outside coverage, pay the distance to the nearest ball

# Hybrid Clustering

- Interpolates between $k$-Median and $k$-Center

- Think of placing $k$ WiFi routers, each with coverage radius $r$

- Clients within coverage, pay 0 (zero)

- Clients outside coverage, pay the distance to the nearest ball



$$d_r(p, X) := max\{d(p, X) - r, 0\} \qquad \text{for } X \subseteq F$$

# Hybrid Clustering

- $d_r(p, X) := max\{d(p, X) - r, 0\}$   $r$-distance

- Input

  $P$: set of $n$ points

  $F$: set of facilities

  $k$: positive integer

  $d$: distance function on $P \cup F$

  $r$: non-negative real

- Output
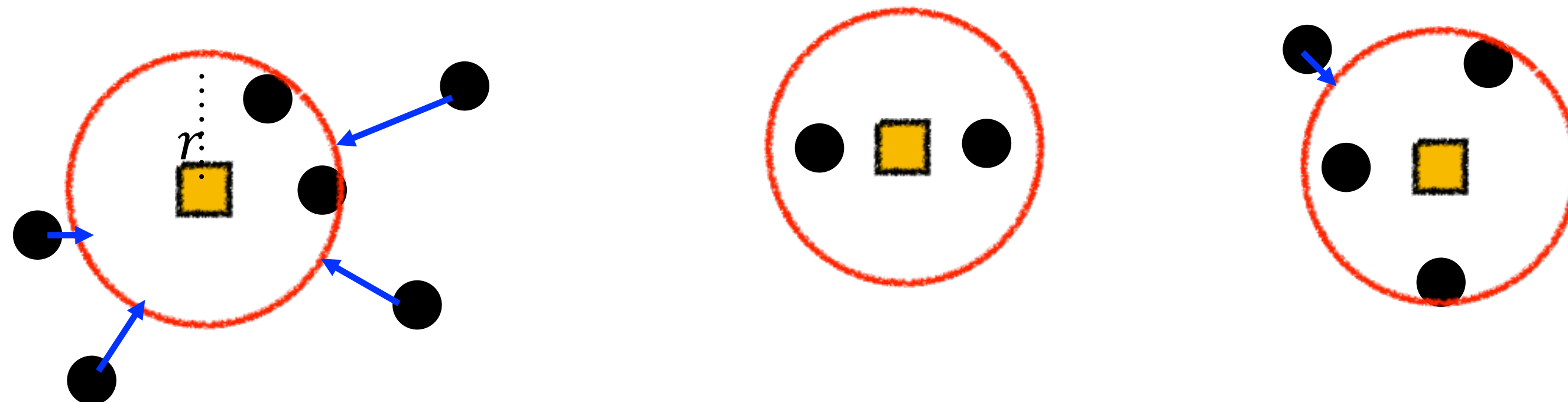
  $X \subseteq F$: set of $k$ centers

- Minimize

  $$\sum_{p \in P} d_r(p, X)$$

# Hybrid Clustering

**Motivation**

- Interpolates between $k$-Median and $k$-Center

- Shape Fitting

  - Extension of Linear regression: Fitting "best" lines

  - Projective Clustering: Fitting "best" affine spaces

# Hybrid Clustering

**Motivation**

- Interpolates between $k$-Median and $k$-Center

- Shape Fitting

  - Extension of Linear regression: Fitting "best" lines

  - Projective Clustering: Fitting "best" affine spaces

  - Hybrid Clustering: Fitting "best" $r$-radius balls

# Hybrid Clustering

**Motivation**

- <span style="color:#e0218a">Interpolates</span> between $k$-Median and $k$-Center

- <span style="color:#e0218a">Shape Fitting</span>
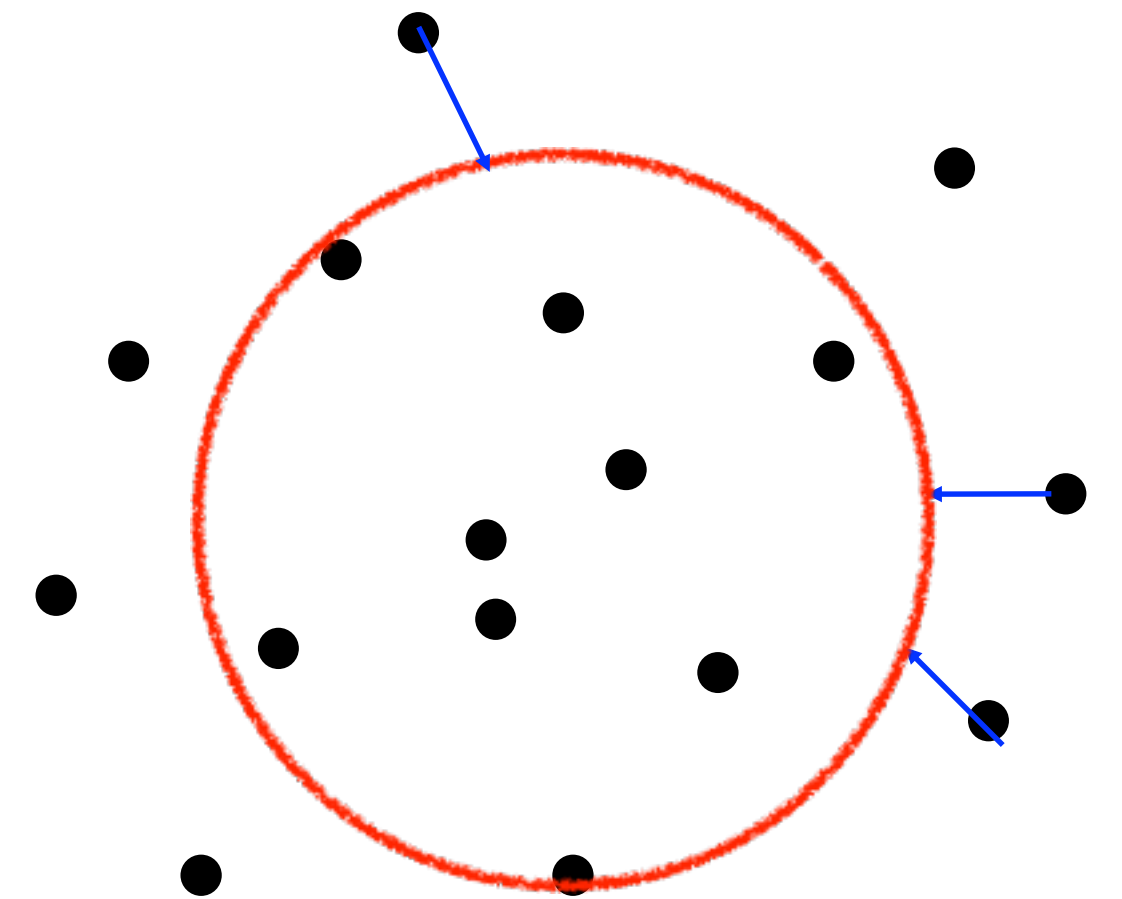
  - Extension of <span style="color:blue">Linear regression:</span> Fitting "best" lines

  - <span style="color:blue">Projective Clustering:</span> Fitting "best" affine spaces

  - <span style="color:blue">Hybrid Clustering</span>: Fitting "best" $r$-radius balls

# Literature

- Recently introduced by Fomin, Golovach, Inamdar, Saurabh, Zehavi [Approx' 24]

- $r = 0: d_r(p, X) = d(p, X) \implies k\text{-Median}$

- $r = OPT_{kc}: \sum d_r(p, X) = 0 \implies k\text{-Center}$    $OPT_{kc} = k\text{-Center OPT}$

- No Uni-criteria approximations:    have to violate both—cost & radius

*Hybrid k-Clustering: Blending k-Median and k-Center*

# Literature

- Recently introduced by Fomin, Golovach, Inamdar, Saurabh, Zehavi [Approx' 24]

- No Uni-criteria approximations:    have to violate both—cost & radius

# Literature

- Recently introduced by Fomin, Golovach, Inamdar, Saurabh, Zehavi [Approx' 24]

- No Uni-criteria approximations:    have to violate both—cost & radius

- Studied the problem in $\mathbb{R}^d$, where centers can be chosen anywhere

- For $\mathbb{R}^d$, designed $(1 + \epsilon, 1 + \epsilon)$-bicritera approximation

  $\boxed{OPT \text{ cost using } r\text{-radius balls}}$

  - whose cost using $(1 + \epsilon)r$-radius balls is at most $(1 + \epsilon)\text{OPT}_r$

  - in time $\text{FPT}(k, d, \epsilon)$

# Our Results

Theorem 1.

Substantially improve and generalize the results of Fomin at al.

# Our Results

no $d$ here

**Theorem 1.**

For $\mathbb{R}^d$, design $(1 + \epsilon, 1 + \epsilon)$-bicritera approximation in time $FPT(k, \epsilon)$

**Fedor et al [Approx'24].**

For $\mathbb{R}^d$, design $(1 + \epsilon, 1 + \epsilon)$-bicritera approximation in time $2^{(kd/\epsilon)^{O(1)}} n^{O(1)}$

# Our Results

Theorem 1.

For $\mathbb{R}^d$, design $(1+\epsilon, 1+\epsilon)$-bicritera approximation in time $2^{\tilde{O}(k/\epsilon^5)} n^{O(1)}$

no $d$ here

Fedor et al [Approx'24].

For $\mathbb{R}^d$, design $(1+\epsilon, 1+\epsilon)$-bicritera approximation in time $2^{(kd/\epsilon)^{O(1)}} n^{O(1)}$

# Our Results

Theorem 1.

no $d$ here

> For $\mathbb{R}^d$, design $(1 + \epsilon, 1 + \epsilon)$-bicritera approximation in time $FPT(k, \epsilon)$

Works for metric spaces with bounded (algorithmic) scatter dimension

| Bounded Doubling | Bounded Treewidth | Planar | Minor-closed |

Works even when the objective is a monotone norm of $r$-distances

Generalizes the FOCS'23 framework of Abbasi* et al. to $r$-distances

*Parameterized Approximation Schemes for Clustering with General Norm Objectives*
Abbasi, Banerjee, Byrka, Chalermsook, G., Khodamoradi, Marx, Sharma, Spoerhase

# Our Results

**Theorem 1.**

no $d$ here

For $\mathbb{R}^d$, design $(1 + \epsilon, 1 + \epsilon)$-bicritera approximation in time $FPT(k, \epsilon)$

Generalizes the FOCS'23 framework of Abbasi* et al. to $r$-distances

**Theorem 2.**

Design coresets of size $2^{O(d\log(1/\epsilon))} k\log n$ in doubling metrics of dimension $d$

# This talk

Theorem 1.

For $\mathbb{R}^d$, design $(1+\epsilon, 1+\epsilon)$-bicritera approximation in time $FPT(k, \epsilon)$

# This talk

Theorem 1.

For $\mathbb{R}^d$, design $(1+\epsilon, 1+\epsilon)$-bicritera approximation in time $FPT(k, \epsilon)$

- Idea based on EPAS framework of Abbasi et al. [FOCS'23],   Unified-EPAS

  - $(1+\epsilon)$-approximation running in time $FPT(k, \epsilon)$

  - for many clustering problems

  - under any metric space that has bounded (algorithmic) scatter dimension

  - in a unified manner

EPAS: Efficient Parameterized Approximation Schemes
FPT-AS

# Unified-EPAS: Basic Idea

Consider the clustering corresponding to an optimal solution $O$

For each cluster $j \in [k]$, we maintain a cluster constraint $Q_j$

Each $Q_j$ is a sequence of pairs $(p, r_p)$, where $p \in$ Cluster $j$ and $r_p \leq d(p, O)$

$o_1$

$o_2$

$o_3$

# Unified-EPAS: Basic Idea

Consider the clustering corresponding to an optimal solution $O$

For each cluster $j \in [k]$, we maintain a cluster constraint $Q_j$

Each $Q_j$ is a sequence of pairs $(p, r_p)$, where $p \in$ Cluster $j$ and $r_p \leq d(p, O)$

# Unified-EPAS: Basic Idea

Consider the clustering corresponding to an optimal solution $O$

For each cluster $j \in [k]$, we maintain a cluster constraint $Q_j$

Each $Q_j$ is a sequence of pairs $(p, r_p)$, where $p \in$ Cluster $j$ and $r_p \leq d(p, O)$

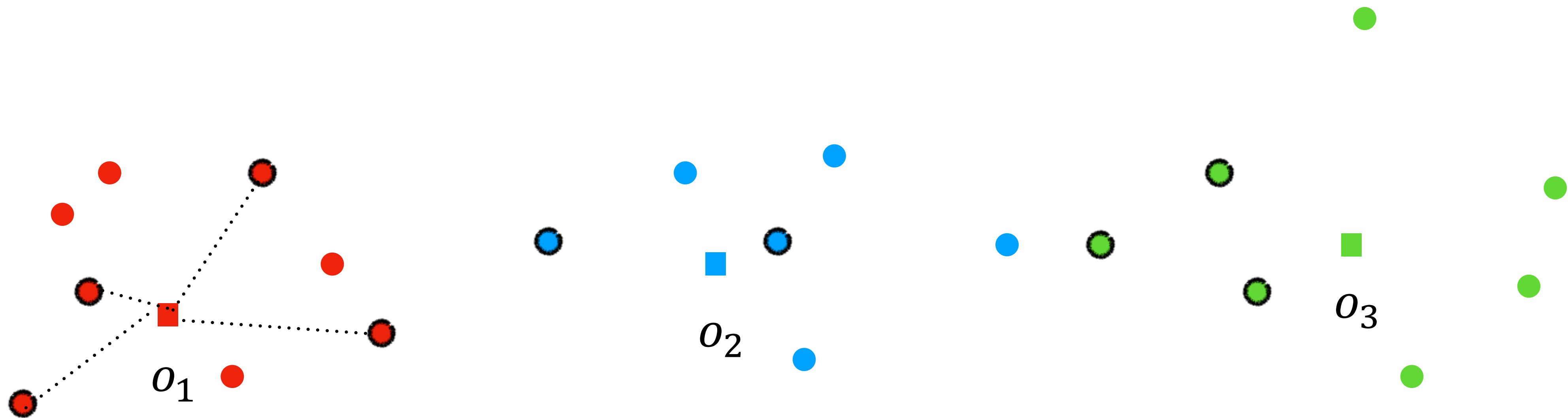Find $X = (x_1, \ldots, x_k)$ such that $x_i$ satisfies all requests in $Q_i$

# Unified-EPAS

# Unified-EPAS

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

$x_1$

$o_1$

$x_2$

$o_2$

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

# Unified-EPAS



$x_1$

$o_1$

$x_2$

$o_2$

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$

**No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

# Unified-EPAS

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ — No → Return $X$

Yes

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

$x_1$

$x_2$

$o_1$

$o_2$

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$      $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$      Update cluster constraint $Q_j$

$x_1$    $x_2$    $o_1$    $o_2$

# Unified-EPAS



Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$     $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

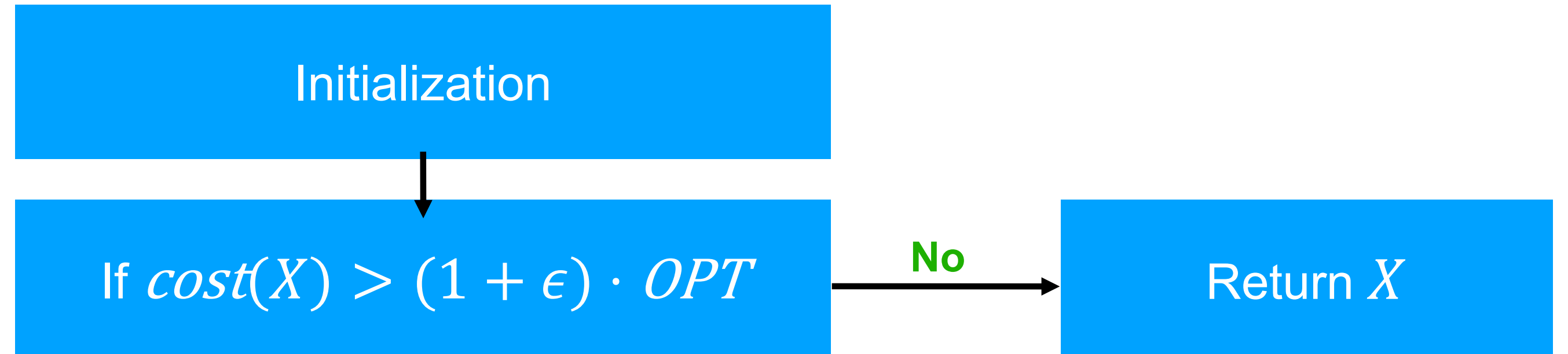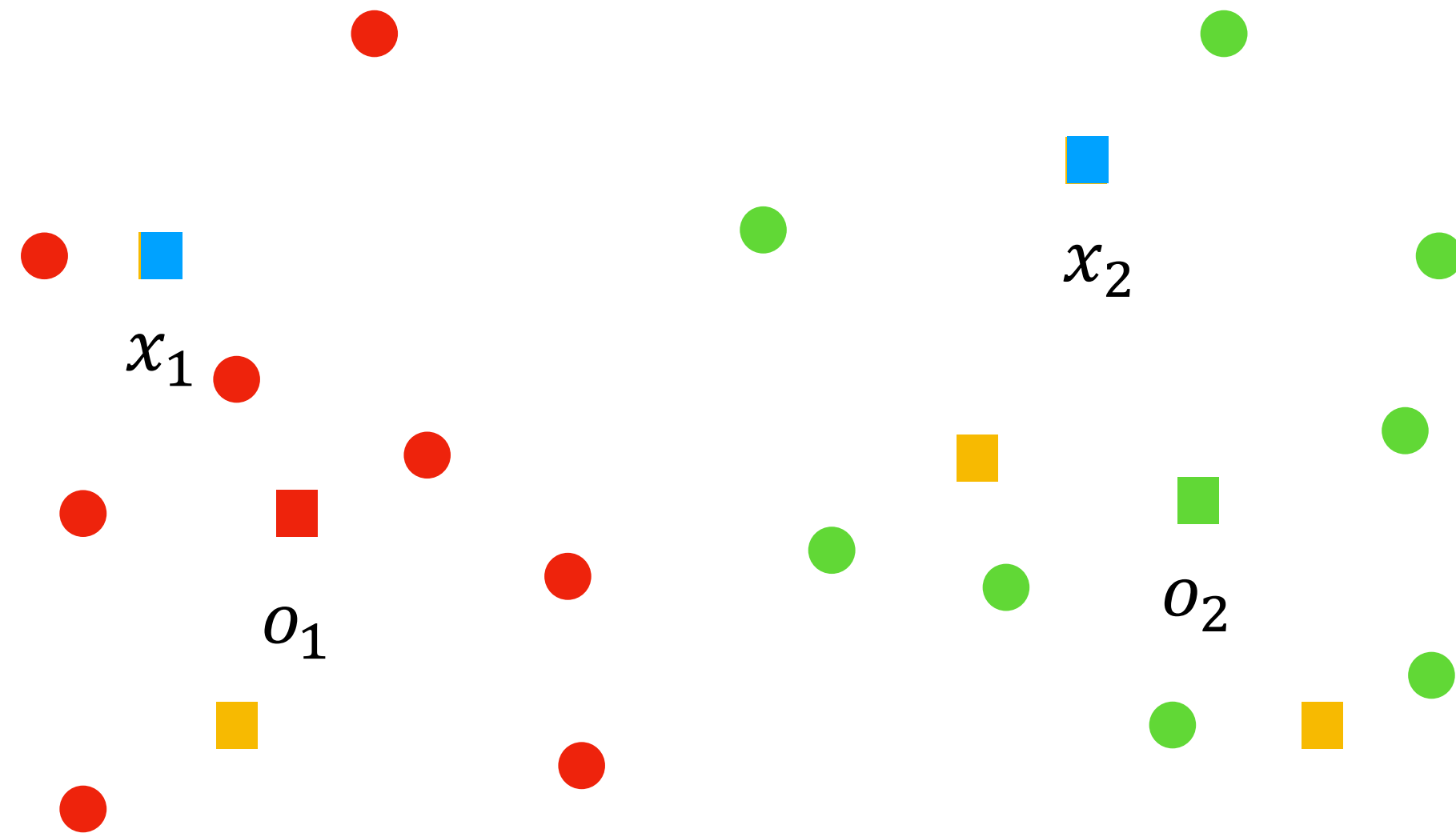Add $\left(p, \frac{d(p,X)}{1 + \epsilon/10}\right)$ to $Q_j$     Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS



Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$     $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$     Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$  — **No** → **Return $X$**

**Yes**

Find a "witness" $p \in P$ to $X$

$$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Update cluster constraint $Q_j$

Recompute $x_j$

$x_2$

$o_1$

$o_2$

$x_1$

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$  →  **No**  →  **Return $X$**

**Yes**

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Update cluster constraint $Q_j$

Recompute $x_j$

$o_1$

$o_2$

$x_2$

$x_1$

# Unified-EPAS



**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$     $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$     Update cluster constraint $Q_j$

Recompute $x_j$

$o_1$   $o_2$   $x_2$   $x_1$

# Unified-EPAS



**Initialization**

**If $cost(X) > (1 + \epsilon) \cdot OPT$** — **No** → **Return $X$**

**Yes**

**Find a "witness" $p \in P$ to $X$**

$$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$$

**Guess cluster $j \in [k]$ of $p$ in $O$**

**Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$**

Update cluster constraint $Q_j$

**Recompute $x_j$**

# Unified-EPAS

**Lemma 1**

If $cost(X) > (1 + \epsilon) \cdot OPT$, then we can find a witness to $X$ w.h.p.

**Question:**

Bound #iterations?

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left( p, \frac{d(p,X)}{1+\epsilon/10} \right)$ to $Q_j$

Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS

If $cost(X) > (1 + \epsilon) \cdot OPT$, then we can find a witness to $X$ w.h.p.

**Question:**

Bound #iterations?

$\epsilon$-scatter dimension

Upper Bounds

**Initialization**

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS

Bound #iterations?

$\epsilon$-scatter dimension

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ → **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$        $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left( p, \frac{d(p,X)}{1+\epsilon/10} \right)$ to $Q_j$        Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS

# Unified-EPAS

Bound #iterations?

$\epsilon$-scatter dimension

Fix $Q_j$

$x_j^1$ ▪ ——— $> (1 + \epsilon)r_p^1$ ● $(p_j^1, r_p^1)$

$x_j^2$ ▪ ——— $\le r_p^1$ ● $(p_j^2, r_p^2)$

$x_j^3$ ▪ ———————— ● $(p_j^2, r_p^3)$

$x_j^4$ ▪ ——— $(1 + \epsilon)r_p^4$ ● $(p_j^4, r_p^4)$

$x_j^5$ ▪ ———————— ● $(p_j^5, r_p^5)$

$x_j^6$ ▪ ——— $(1 + \epsilon)r_p^6$ ● $(p_j^6, r_p^6)$

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ —— **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$    $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1 + \epsilon/10}\right)$ to $Q_j$    Update cluster constraint $Q_j$

Recompute $x_j$

# Unified-EPAS

Fix $Q_j$

$x_j^1$ ■ —————— $> (1+\epsilon)r_p^1$ ● $(p_j^1, r_p^1)$

$x_j^2$ ■ ——— $\leq r_p^1$ ———— ● $(p_j^2, r_p^2)$

$x_j^3$ ■ ———————————— ● $(p_j^2, r_p^3)$

$x_j^4$ ■ —— $(1+\epsilon)r_p^4$ —— ● $(p_j^4, r_p^4)$

$x_j^5$ ■ ———————————— ● $(p_j^5, r_p^5)$

$x_j^6$ ■ —— $(1+\epsilon)r_p^6$ —— ● $(p_j^6, r_p^6)$

$\epsilon$-scattering

**Initialization**

If $cost(X) > (1+\epsilon) \cdot OPT$ —— **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$ $\qquad d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$ $\qquad$ Update cluster constraint $Q_j$

Recompute $x_j$

$\epsilon$-scatter dimension of a metric space is $\lambda$ $\qquad$ if any $\epsilon$-scattering contains at most $\lambda$ many triples with same radius

# Unified-EPAS

$\epsilon$-scatter dimension

Fix $Q_j$

$x_j^1$ ▪ $> (1+\epsilon)r_p^1$ ● $(p_j^1, r_p^1)$

$x_j^2$ ▪ $\leq r_p^1$ ● $(p_j^2, r_p^2)$

$x_j^3$ ▪ ● $(p_j^2, r_p^3)$

$x_j^4$ ▪ $(1+\epsilon)r_p^4$ ● $(p_j^4, r_p^4)$

$x_j^5$ ▪ ● $(p_j^5, r_p^5)$

$x_j^6$ ▪ $(1+\epsilon)r_p^6$ ● $(p_j^6, r_p^6)$

$\epsilon$-scattering

Initialization

If $cost(X) > (1+\epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$ $\qquad$ $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \dfrac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$ $\qquad$ Update cluster constraint $Q_j$

Recompute $x_j$

$\epsilon$-scatter dimension of a metric space is $\lambda$ ⇨ if radius aspect ratio is bounded, then the length is bounded

# Unified-EPAS

Bound #iterations?

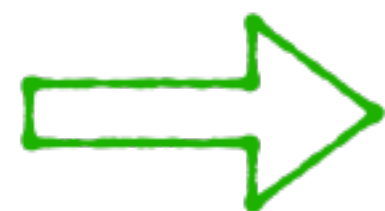Upper Bounds

**Initialization**

Compute Upper bounds

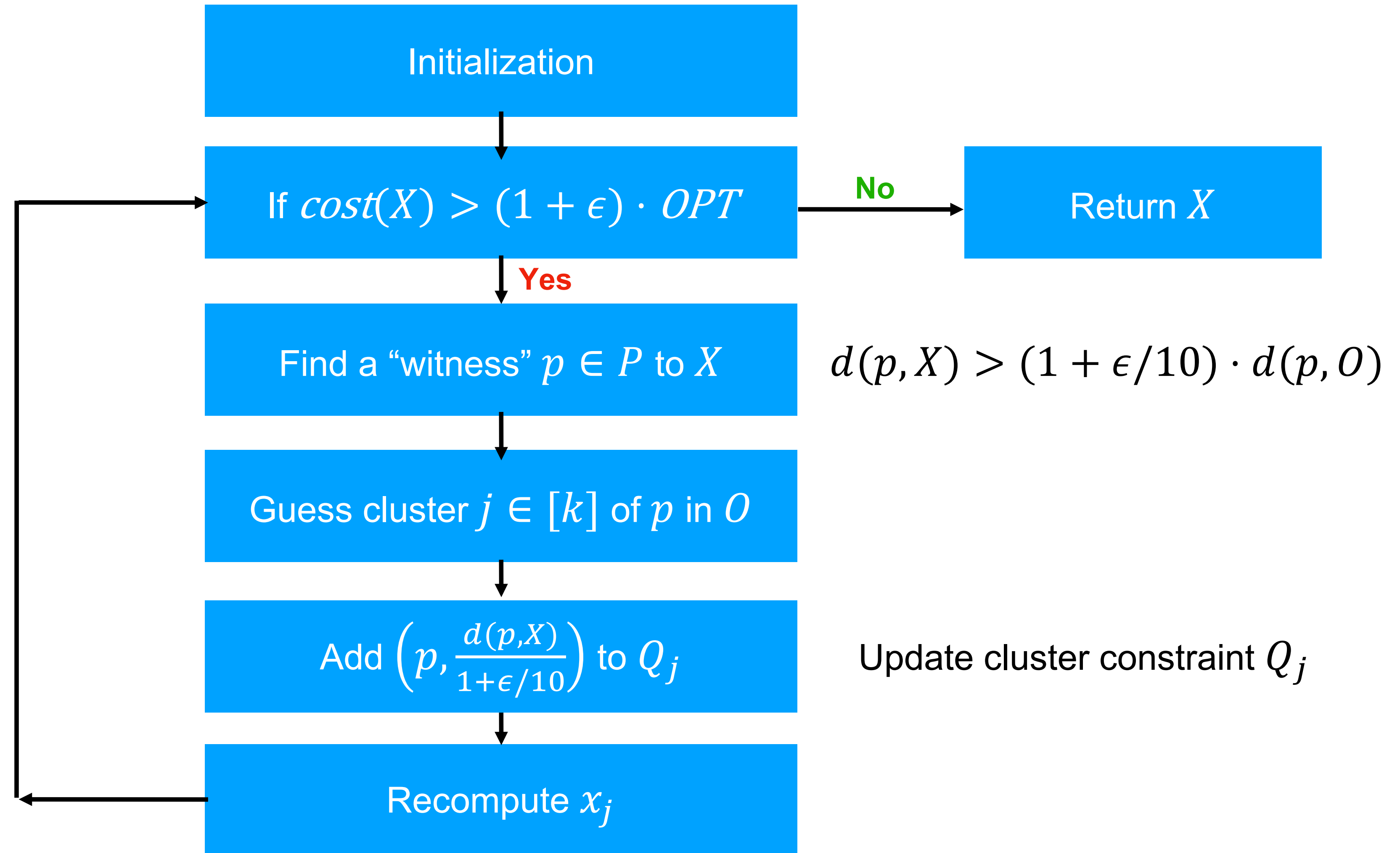Initialize Cluster constraints $Q_1, \ldots, Q_k$ using upper bounds

Initialize solution $X = (x_1, \ldots, x_k)$ using $Q_1, \ldots, Q_k$

**Lemma 2**

Upper bounds

Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ → **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left( p, \frac{d(p,X)}{1+\epsilon/10} \right)$ to $Q_j$

Update cluster constraint $Q_j$

Recompute $x_j$

Radii aspect ratio of requests in every $Q_j$ is bounded

# Unified-EPAS

**Lemma 1**

If $cost(X) > (1 + \epsilon) \cdot OPT$, then we can find a witness to $X$ w.h.p. $\quad g(k, \epsilon)$
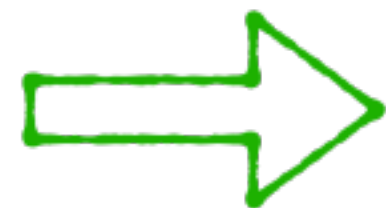
**Lemma 2**

Upper bounds $\Rightarrow$ Radii aspect ratio of requests in every $Q_j$ is bounded $\quad f(k, \epsilon)$

**Lemma 3 (Theorem)**

Lemma 1 ➕ Lemma 2 $\Rightarrow$ Requests in every $Q_j$ form an $\epsilon$-scattering $\quad \lambda(\epsilon)$
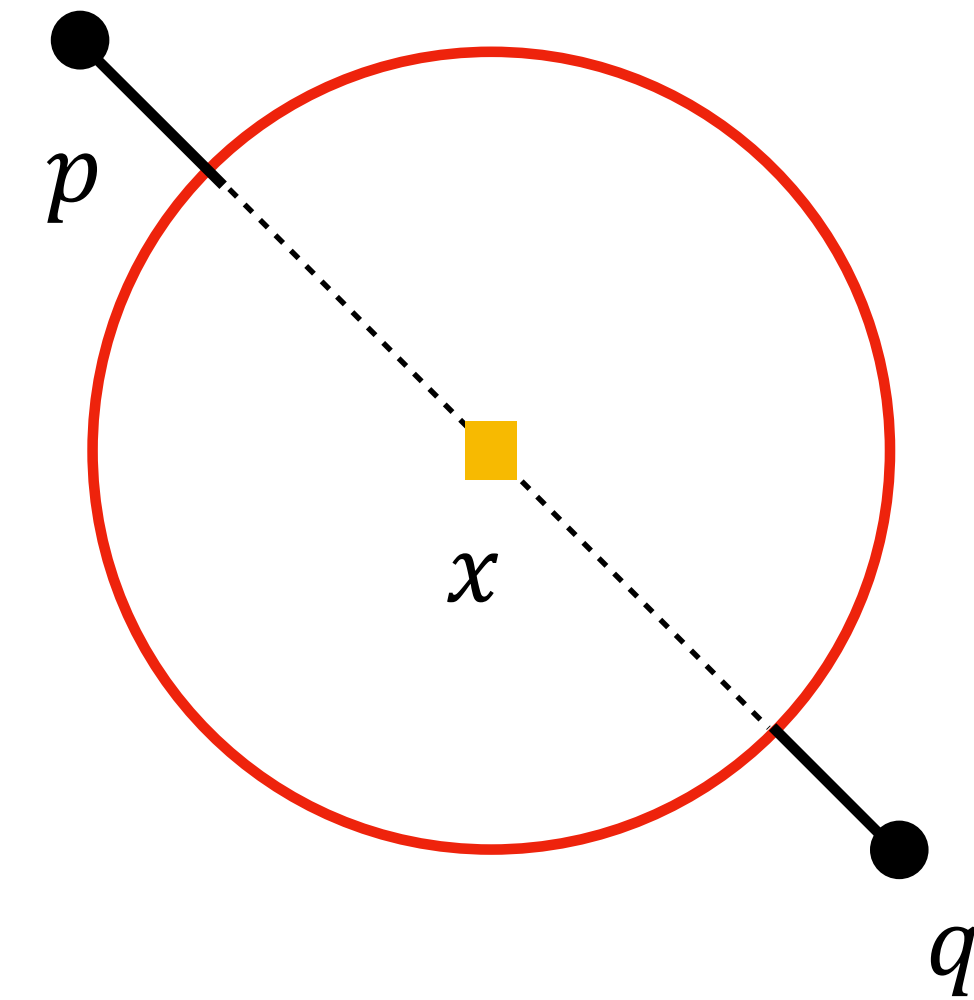
$\Rightarrow$ #iterations are bounded $\quad h(k, \epsilon, \lambda)$

# Hybrid Clustering

$d_r$ does not satisfy triangle inequality

- Computing Upper bounds fails!

- Sampling lemma (Lemma 1) does not work!

- Radii Aspect Ratio lemma (Lemma 2) fails!

- Iteration lemma (Lemma 3) does not apply since the new requests may not be feasible!
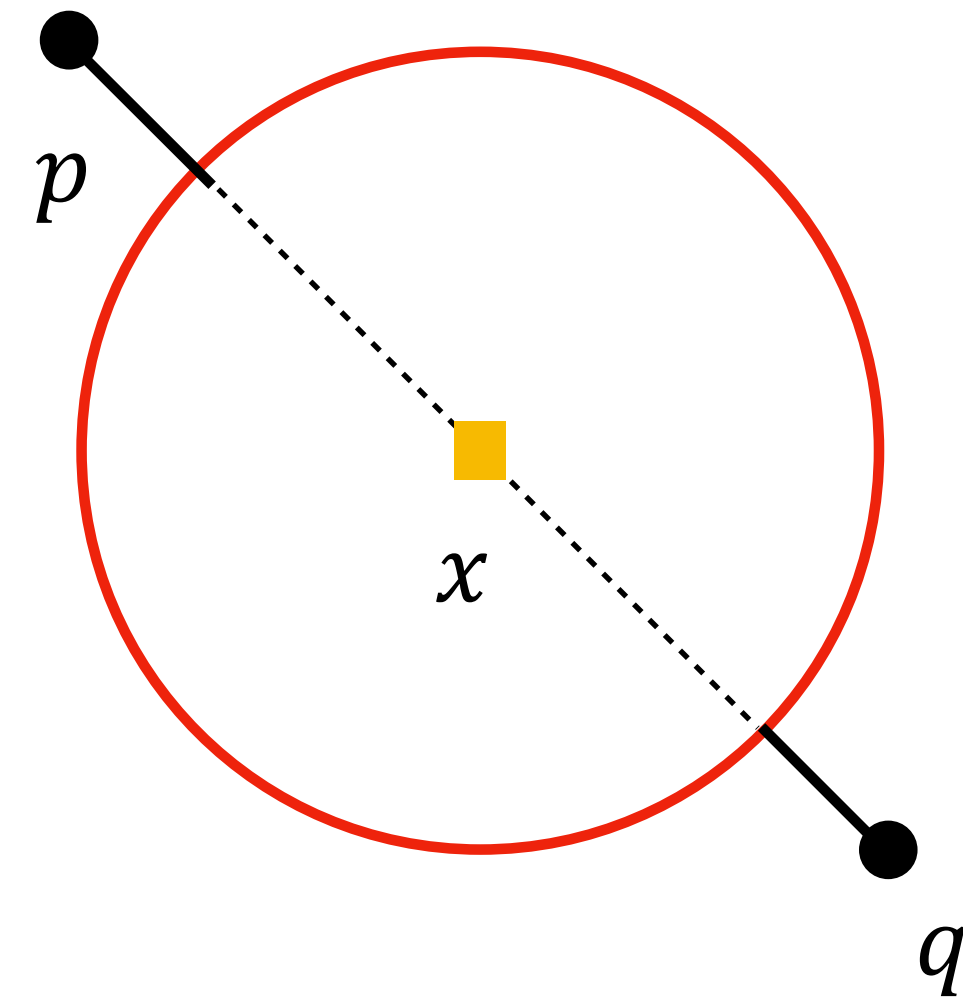
# Hybrid Clustering

$d_r$ does not satisfy triangle inequality

- Computing Upper bounds fails!

- Sampling lemma (Lemma 1) does not work!

- Radii Aspect Ratio lemma (Lemma 2) fails!

- Iteration lemma (Lemma 3) does not apply since the new requests may not be feasible!
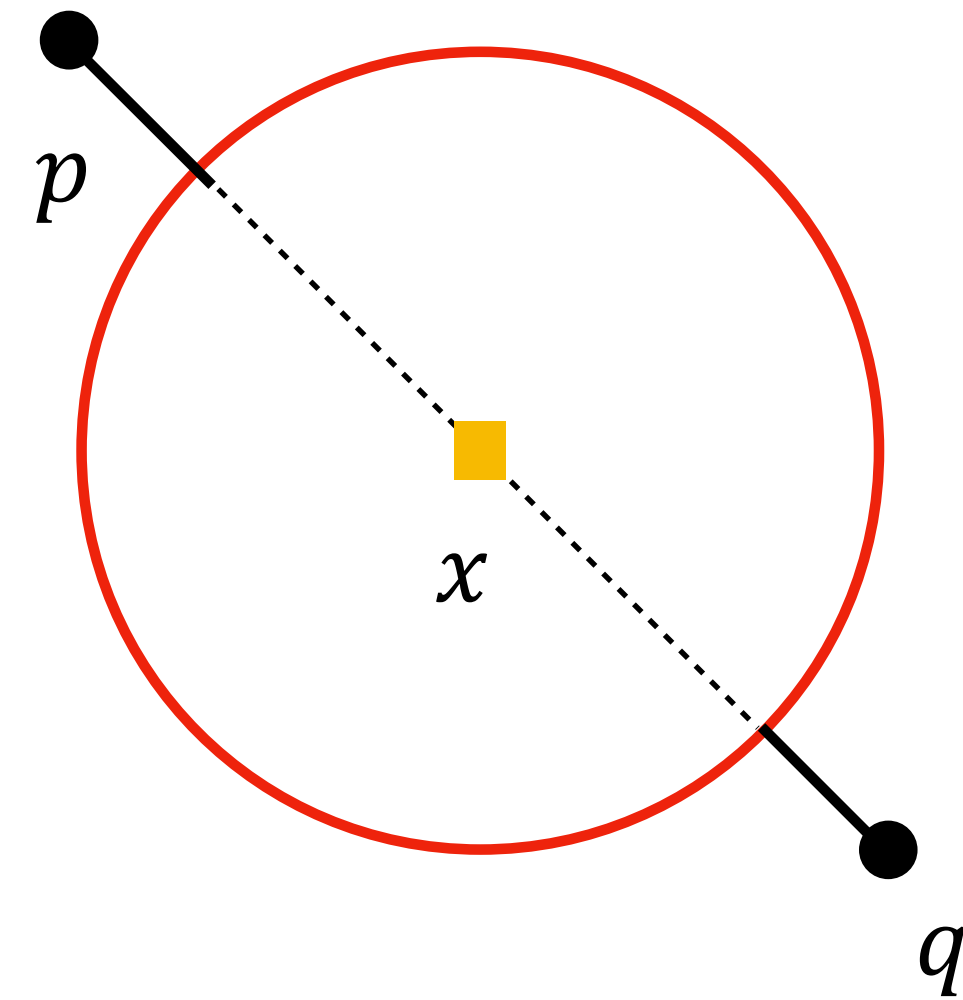
# Hybrid Clustering

$d_r$ does not satisfy triangle inequality



- Computing Upper bounds fails!

- Sampling lemma (Lemma 1) does not work!

- Radii Aspect Ratio lemma (Lemma 2) fails!

- Iteration lemma (Lemma 3) does not apply since the new requests may not be feasible!

# Hybrid Clustering

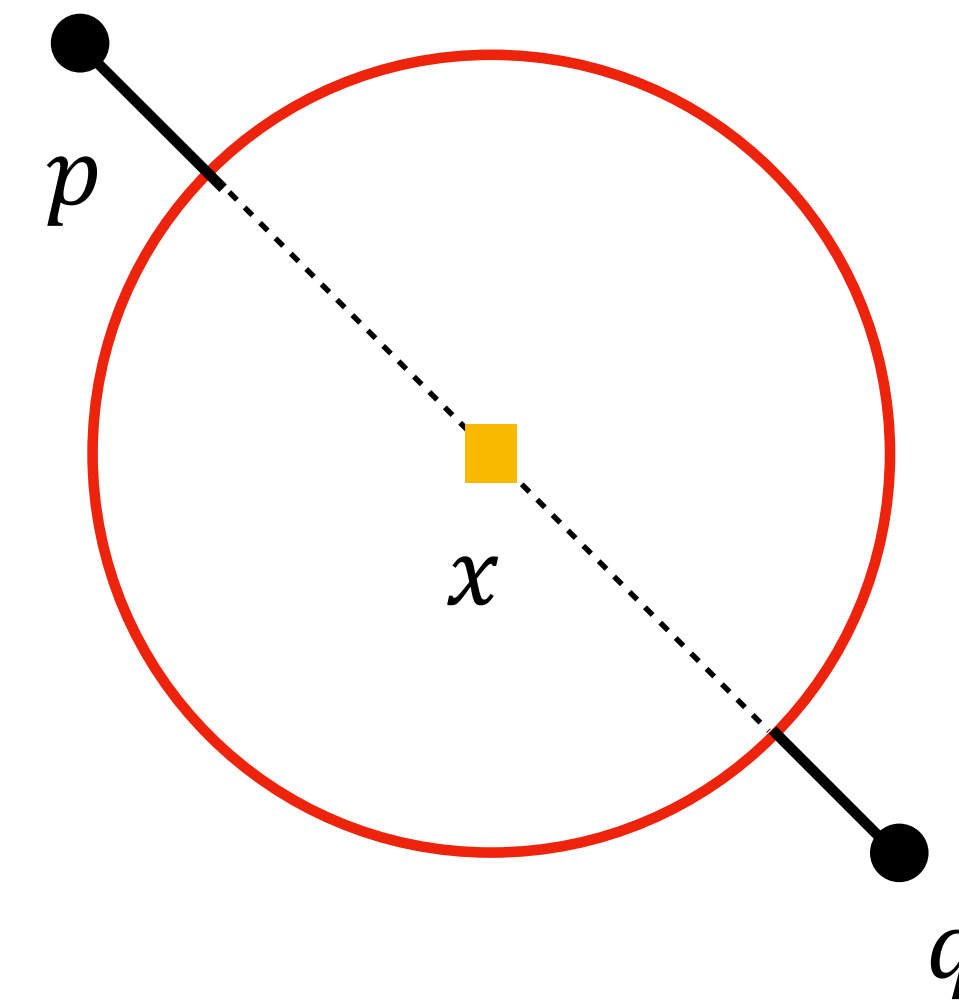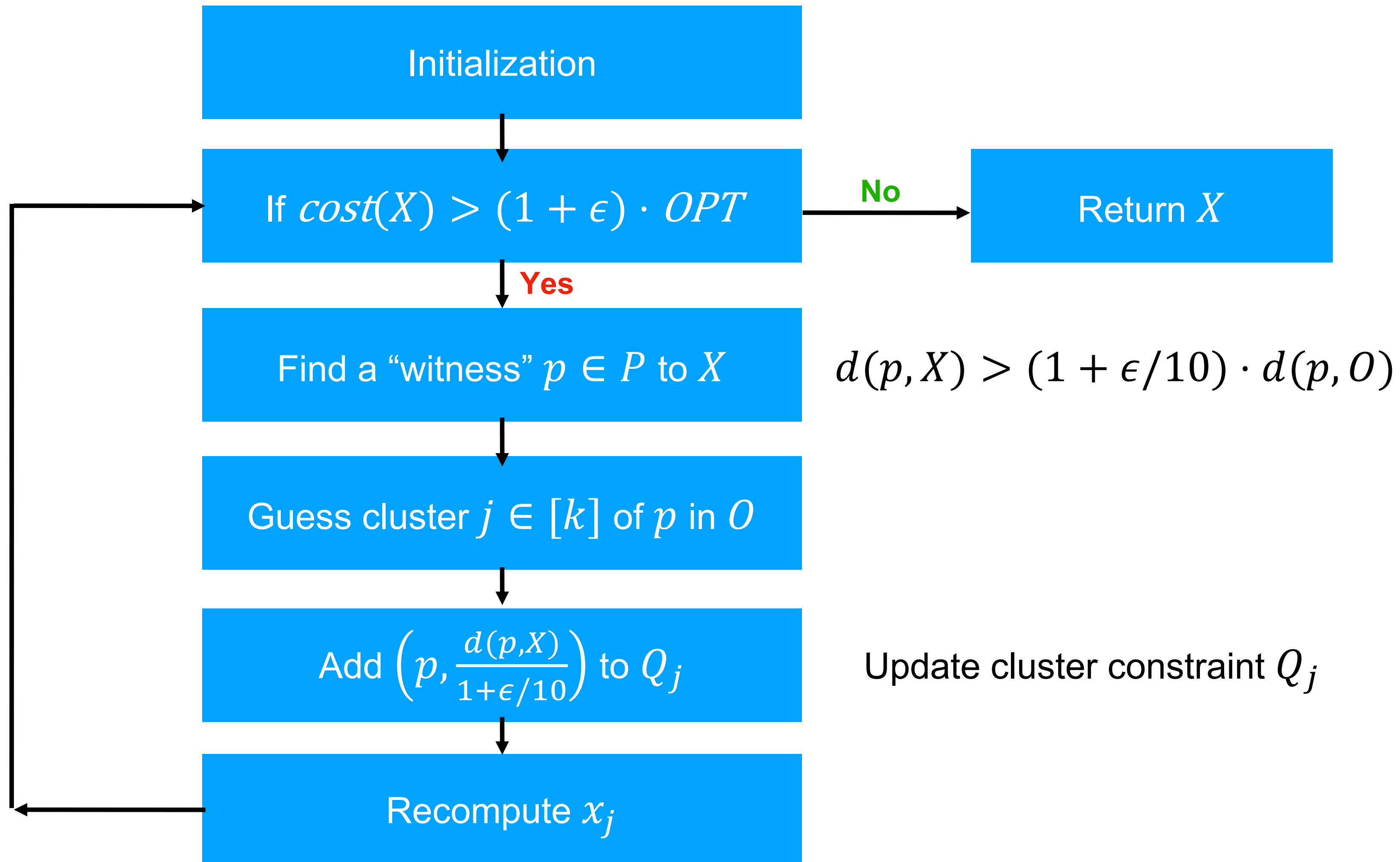$d_r$ does not satisfy triangle inequality

- Computing Upper bounds fails!

- Sampling lemma (Lemma 1) does not work!

- Radii Aspect Ratio lemma (Lemma 2) fails!

- Iteration lemma (Lemma 3) does not apply since the new requests may not be feasible!

# Unified-EPAS



Initialization

If $cost(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$     $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$     Update cluster constraint $Q_j$

Recompute $x_j$

# Attempt 1



Initialization

If $cost_r(X) > (1 + \epsilon) \cdot OPT$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$     $d(p, X) > (1 + \epsilon/10) \cdot d(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$     Update cluster constraint $Q_j$

Recompute $x_j$

# Attempt 1



# Attempt 1

```
Initialization
        │
        ▼
If $cost_r(X) > (1 + \epsilon) \cdot OPT$  ──No──►  Return $X$
        │
       Yes
        │
        ▼
Find a "witness" $p \in P$ to $X$        $d_r(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$
        │
        ▼
Guess cluster $j \in [k]$ of $p$ in $O$
        │
        ▼
Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$        Update cluster constraint $Q_j$
        │
        ▼
Recompute $x_j$
```

# Attempt 1

# Attempt 1



**Initialization**

If $cost_r(X) > (1 + \epsilon) \cdot OPT$ —**No**→ Return $X$

**Yes** ↓

Find a "witness" $p \in P$ to $X$    $d_r(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \dfrac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$    Request may not be feasible!
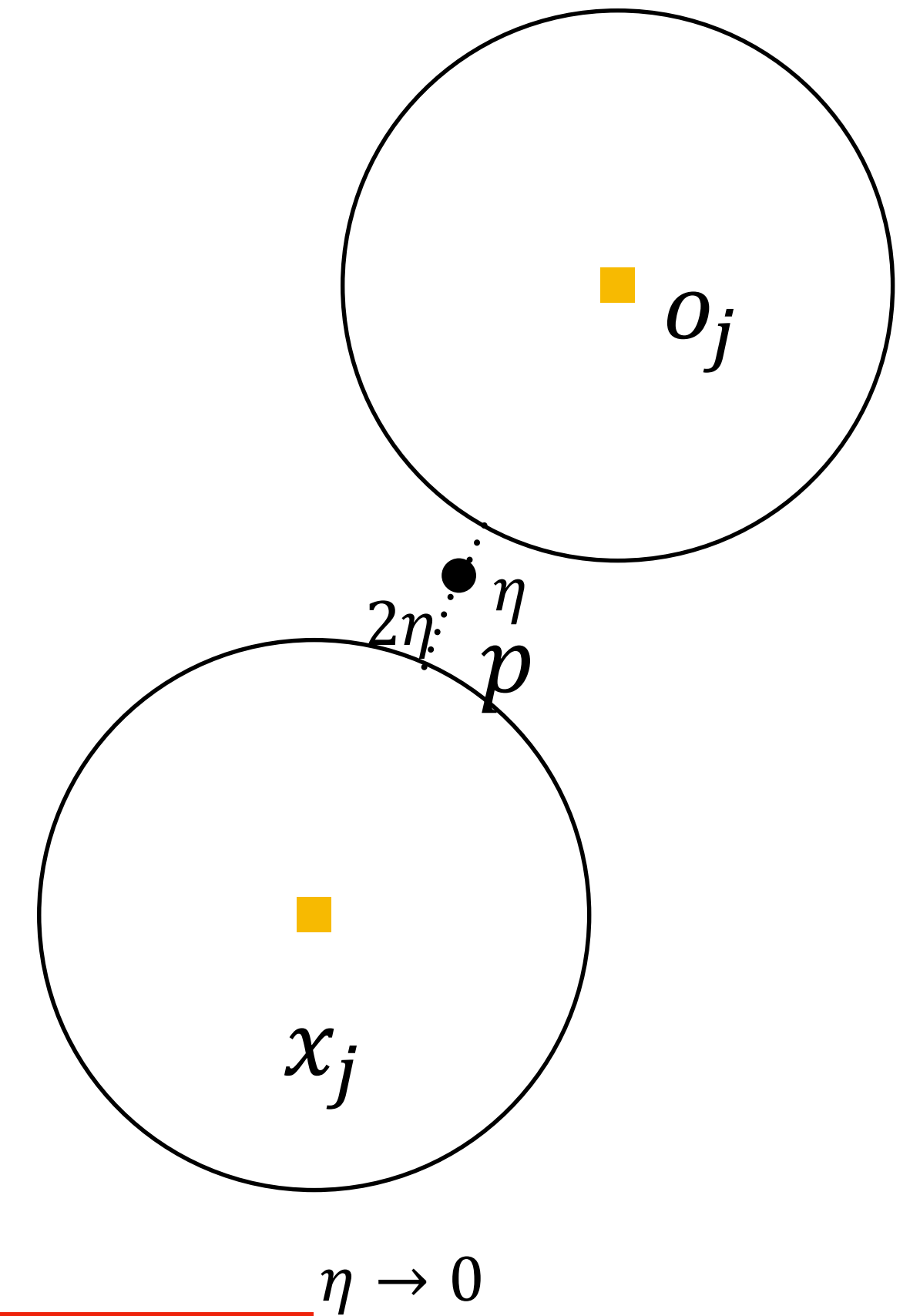
think when $d(p, X) \approx r$

Recompute $x_j$

$o_j$

$2\eta$ $\quad$ $\eta$

$p$

$x_j$

$\eta \to 0$

# Attempt 1



Initialization

If $cost_r(X) > (1 + \epsilon) \cdot OPT$

**No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d_r(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$
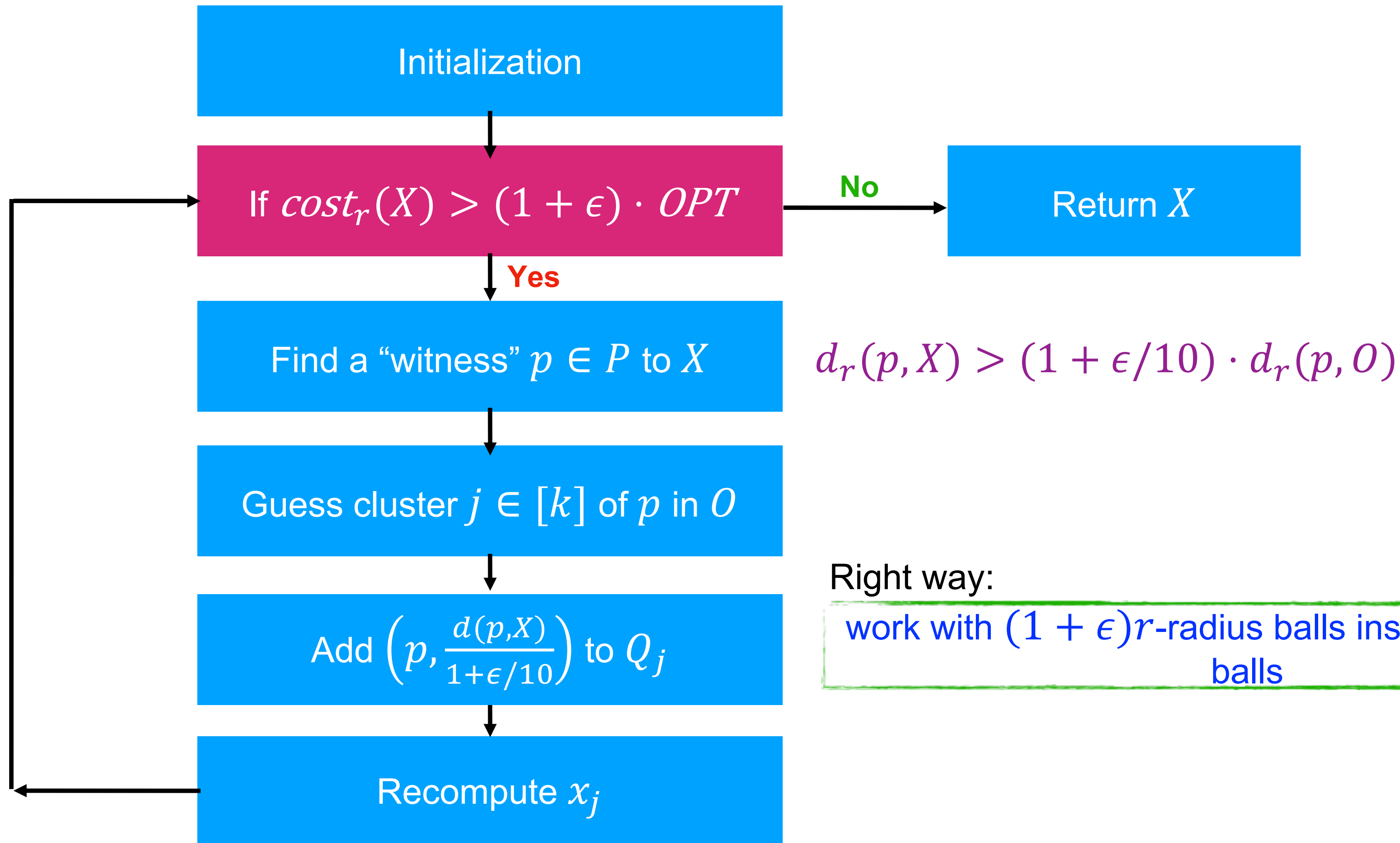
Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Actually, this is more of a fundamental bottleneck

Recompute $x_j$

$\because \implies$ uni-criteria approximation

$o_j$

$2\eta$ $\eta$

$p$

$x_j$

$\eta \to 0$

# Attempt 1

**Initialization**

If $cost_r(X) > (1 + \epsilon) \cdot OPT$ —**No**→ Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$d_r(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \dfrac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Right way:

work with $(1 + \epsilon)r$-radius balls instead of $r$-radius balls

Recompute $x_j$

$o_j$

$2\eta$ $\eta$

$p$

$x_j$

$\eta \to 0$

# Attempt 1

Initialization

If $cost_{(1+\epsilon)r}(X) > (1+\epsilon) \cdot OPT_r$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

$$d_r(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Recompute $x_j$

Right way:

work with $(1+\epsilon)r$-radius balls instead of $r$-radius balls

$\eta \to 0$

$o_j$

$2\eta$   $\eta$

$p$

$x_j$

# Attempt 1

**Initialization**

If $cost_{(1+\epsilon)r}(X) > (1+\epsilon) \cdot OPT_r$ — **No** → **Return $X$**

**Yes**

**Find a "witness" $p \in P$ to $X$**

$$d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$$

**Guess cluster $j \in [k]$ of $p$ in $O$**

**Add $\left(p, \dfrac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$**

**Recompute $x_j$**

$o_j$

$\eta$

$2\eta$

$p$

$x_j$

Right way:

work with $(1+\epsilon)r$-radius balls instead of $r$-radius balls

$\eta \to 0$

# Attempt 1



**Initialization**

If $cost_{(1+\epsilon)r}(X) > (1+\epsilon) \cdot OPT_r$  → **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Recompute $x_j$

$d_{(1+\epsilon)r}(p,X) > (1+\epsilon/10) \cdot d_r(p,O)$

$(1+\epsilon)r$  $r$  $o_j$

$\eta$

$2\eta$  $p$

$(1+\epsilon)r$

$r$  $x_j$

$\eta \to 0$

Right way:

work with $(1+\epsilon)r$-radius balls instead of $r$-radius balls

# Attempt 1



**Initialization**

If $cost_{(1+\epsilon)r}(X) > (1+\epsilon) \cdot OPT_r$ — **No** → Return $X$

**Yes**

Find a "witness" $p \in P$ to $X$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/10}\right)$ to $Q_j$

Recompute $x_j$

$d_{(1+\epsilon)r}(p,X) > (1+\epsilon/10) \cdot d_r(p,O)$

**Right way:**
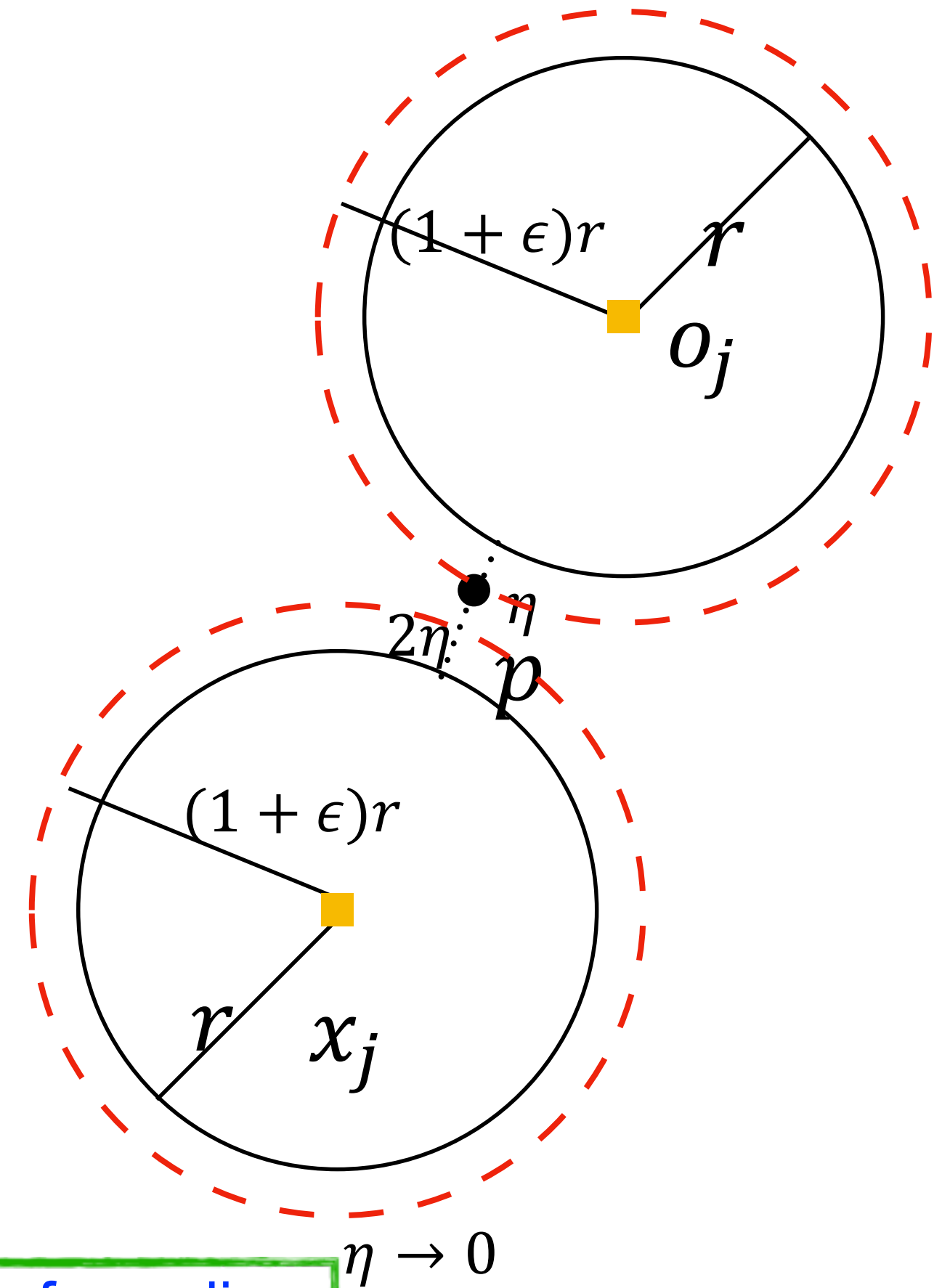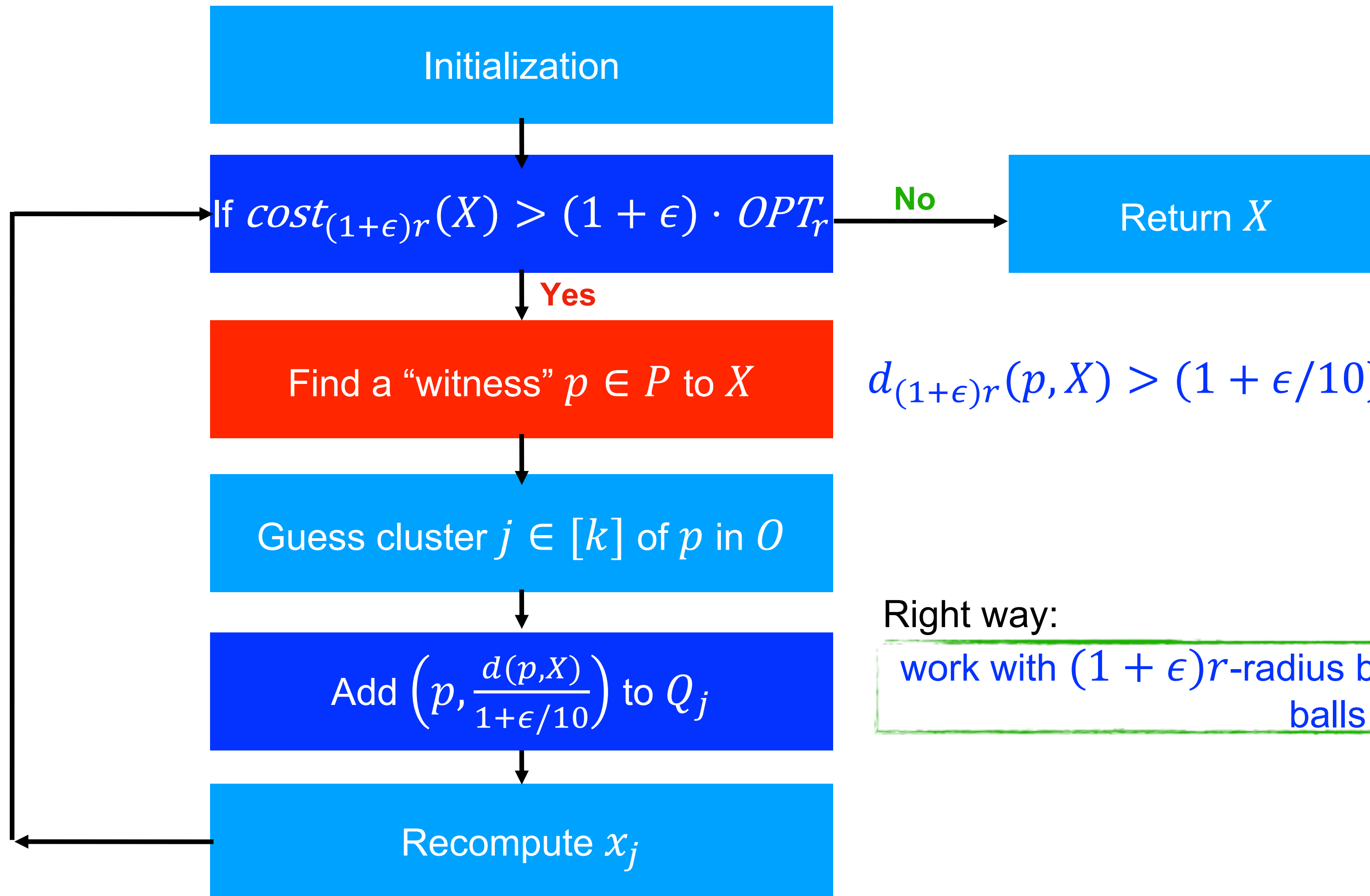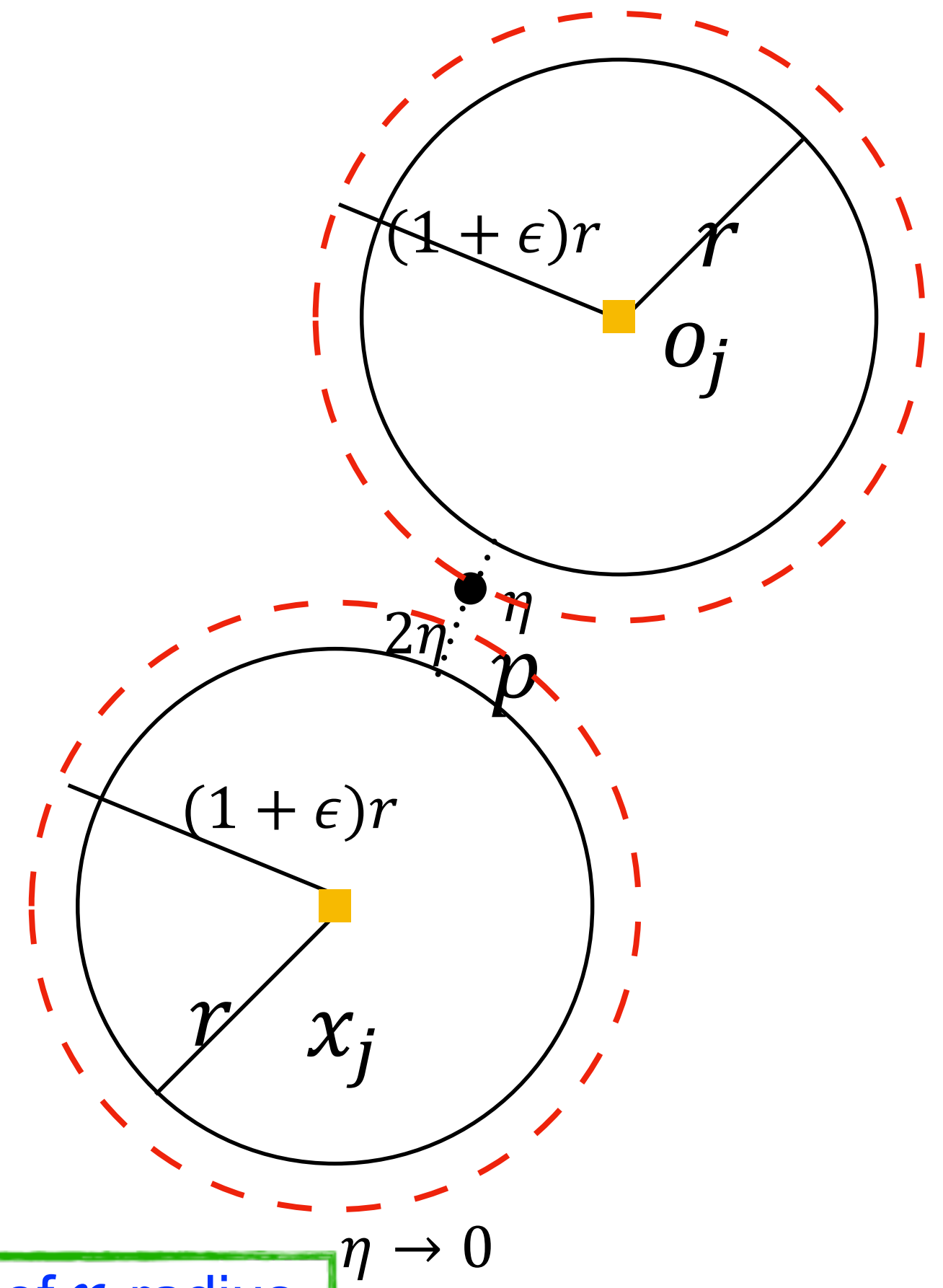work with $(1+\epsilon)r$-radius balls instead of $r$-radius balls

$(1+\epsilon)r$   $r$   $o_j$

$\eta$
$2\eta$   $p$

$(1+\epsilon)r$

$r$   $x_j$

$\eta \to 0$

# Attempt 1

# Sampling Witness

Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

$d_r$ does not satisfy triangle inequality $\implies$ FOCS'23 sampling fails
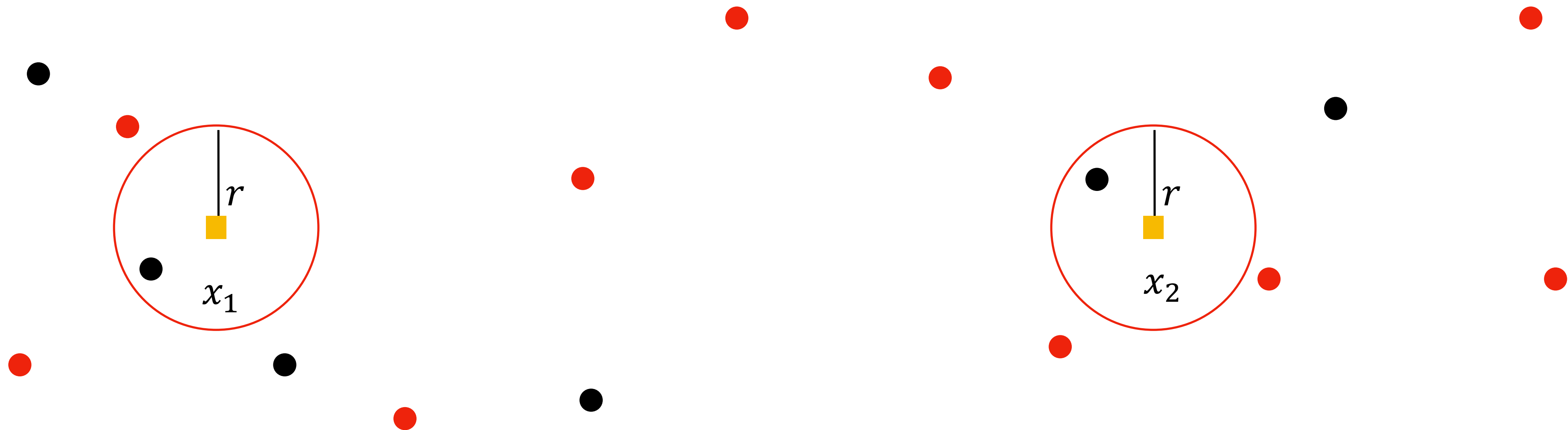
# Sampling Witness

Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$

# Sampling Witness

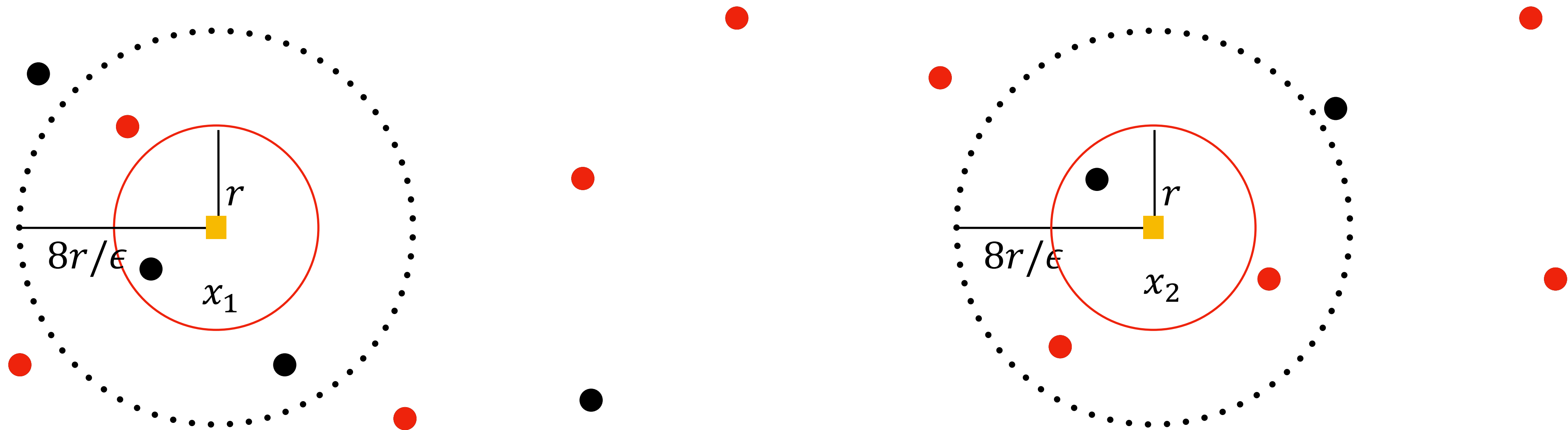Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$

# Sampling Witness

Witness:  $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$

# Sampling Witness

Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$
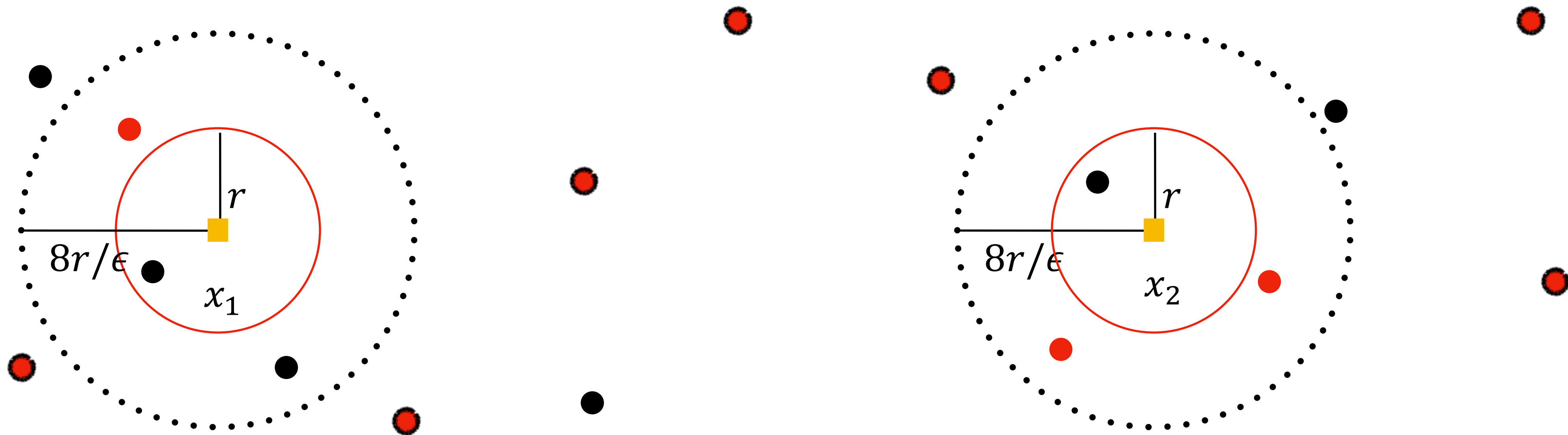
$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$
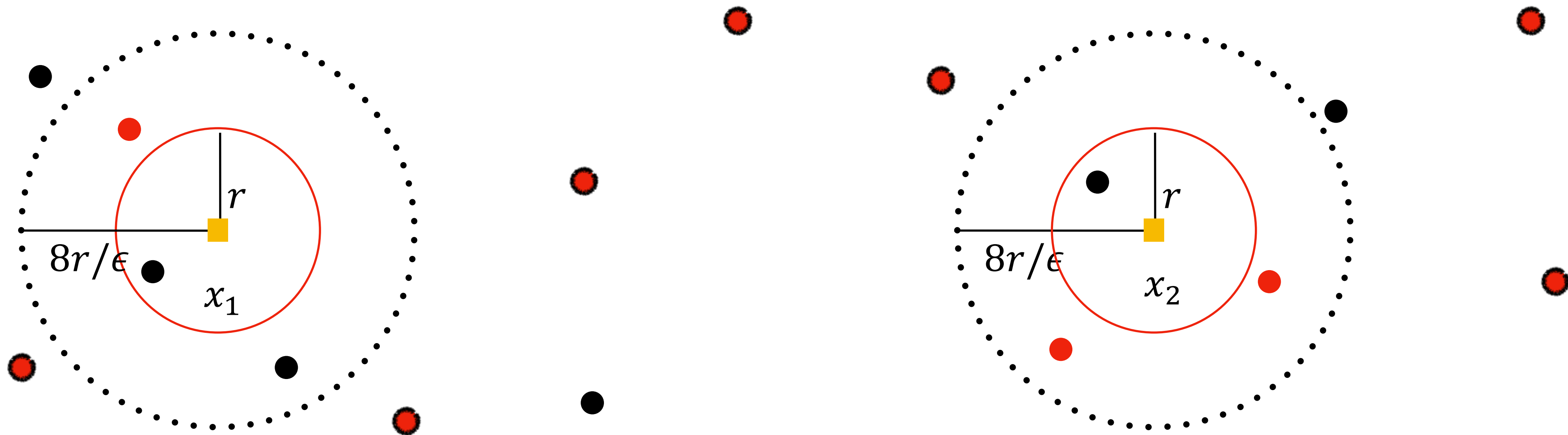
# Sampling Witness

Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

Far away witnesses

$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$
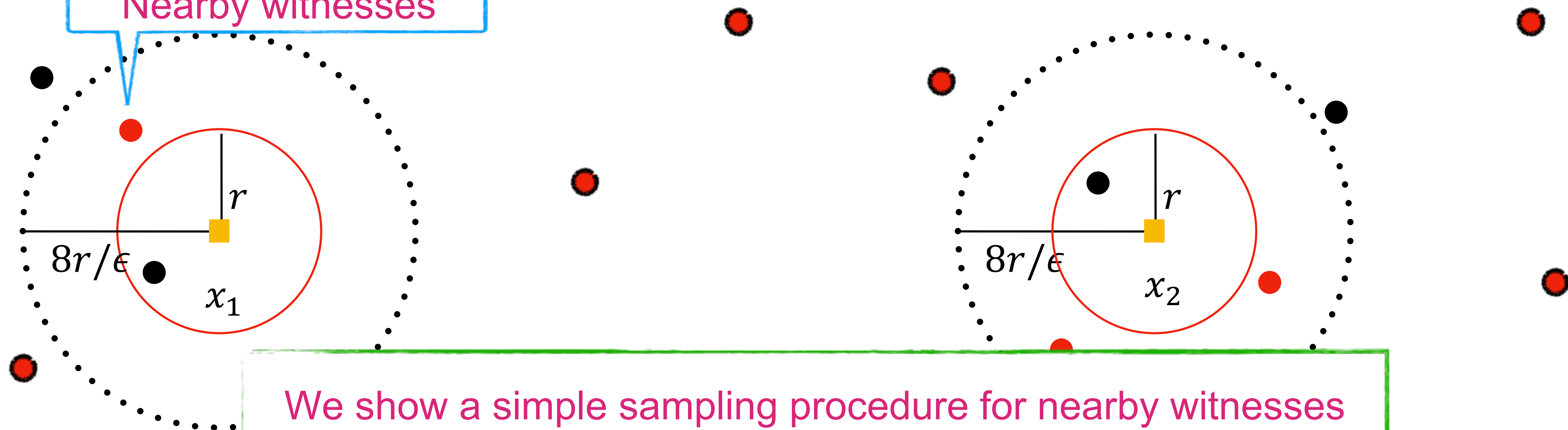
# Sampling Witness

Witness: $d_{(1+\epsilon)r}(p, X) > (1 + \epsilon/10) \cdot d_r(p, O)$

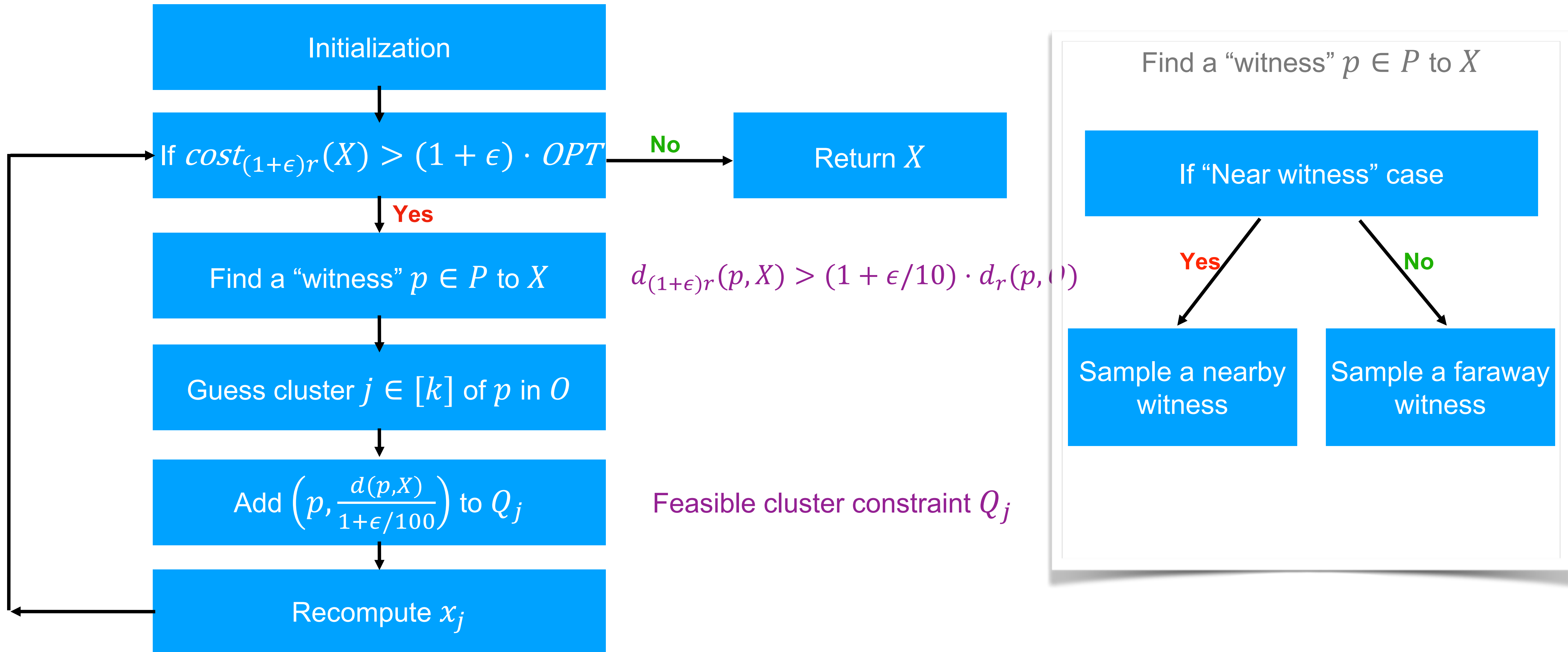$d_r$ does not satisfy triangle inequality. But, $d_r \approx d$ when $d_r = \Omega(r/\epsilon)$

Nearby witnesses

$r$

$8r/\epsilon$

$x_1$

$r$

$8r/\epsilon$

$x_2$

We show a simple sampling procedure for nearby witnesses

# Our Algorithm



Initialization

If $cost_{(1+\epsilon)r}(X) > (1+\epsilon) \cdot OPT$ — No → Return $X$

Yes

Find a "witness" $p \in P$ to $X$

$d_{(1+\epsilon)r}(p,X) > (1+\epsilon/10) \cdot d_r(p,O)$

Guess cluster $j \in [k]$ of $p$ in $O$

Add $\left(p, \frac{d(p,X)}{1+\epsilon/100}\right)$ to $Q_j$

Feasible cluster constraint $Q_j$

Recompute $x_j$

Find a "witness" $p \in P$ to $X$

If "Near witness" case

Yes — Sample a nearby witness

No — Sample a faraway witness

# Summary

Showed a bi-criteria EPAS for Hybrid Clustering

Metric spaces with bounded scatter dimension

Norm objective of $r$-distances

Generalize FOCS'23 EPAS framework for $r$-distances

Designed coresets for Hybrid Clustering in doubling dimensions

Derandomization?

Constrained variants of Hybrid Clustering?

capacities, outlier, fairness

Polynomial-time approximability?

$(18,6)$ is known

*Thank You!*