# Sesión 14.1: Locality-sensitive hashing

**CS3102 EDA**
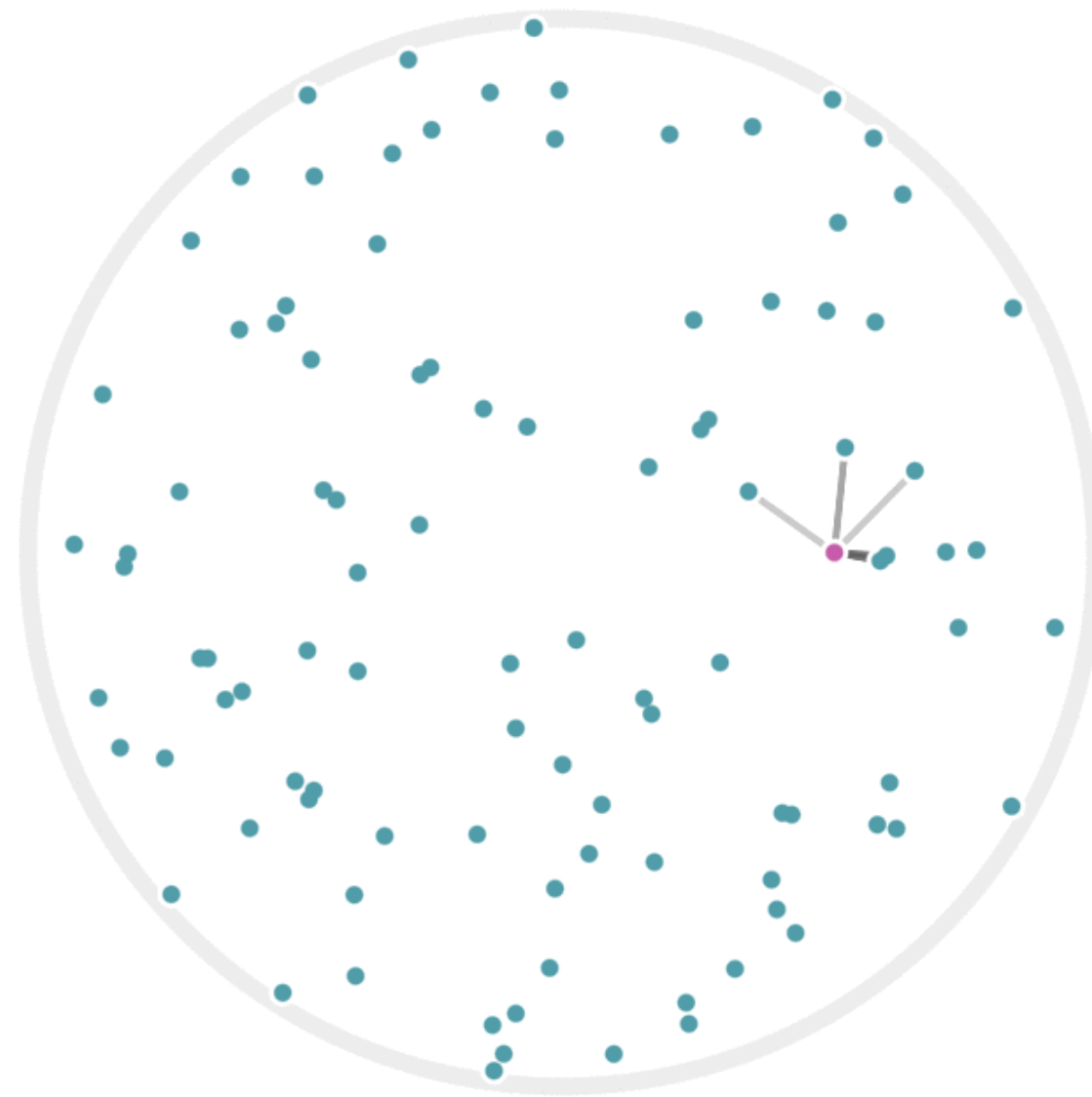
UTEC
UNIVERSIDAD DE INGENIERÍA
Y TECNOLOGÍA

# Índice

Universidad de Ingeniería y Tecnología

# 1. Locality-sensitive hashing (LSH)

# **Locality-sensitive** *hashing (LSH)*



Piotr Indyk and Rajeev Motwani (1998) "Approximate nearest neighbors: towards removing the curse of dimensionality". Proceedings of the thirtieth annual ACM symposium on Theory of computing. p. 604-613.

# **Hash** *Function*

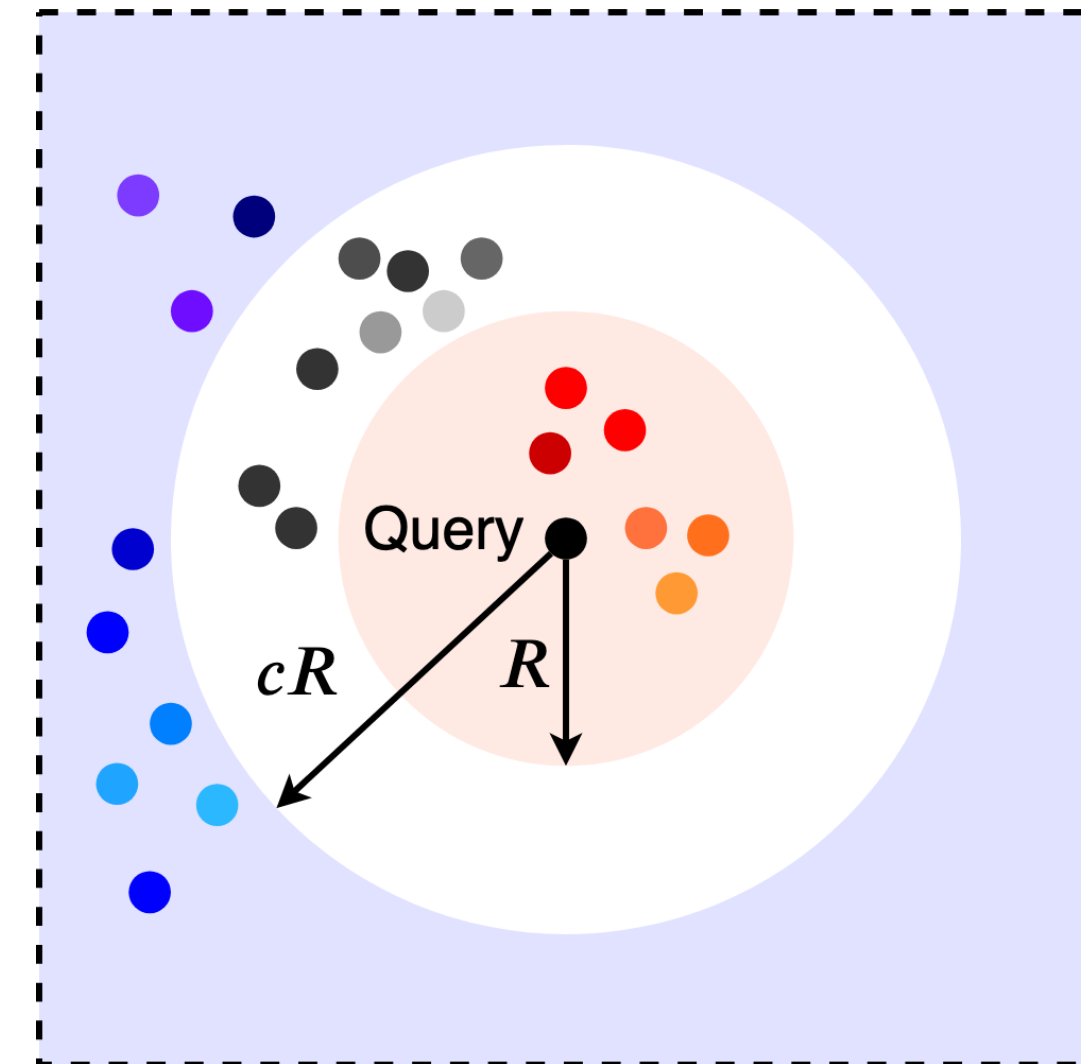# Locality-sensitive *hashing (LSH)*



keys     hashing function     hash buckets     values

# **Locality-sensitive** *hashing (LSH)*

LSH *family* $\mathcal{F}$ es $(R, cR, p_1, p_2)$-sensitive con respecto a la distancia $d(x, y)$ si para algún $h \in \mathcal{H}$ tenemos que:

- Si $d(x, y) \leq R$ entonces $P_{\mathcal{H}}[h(x) = h(y)] \geq p_1$
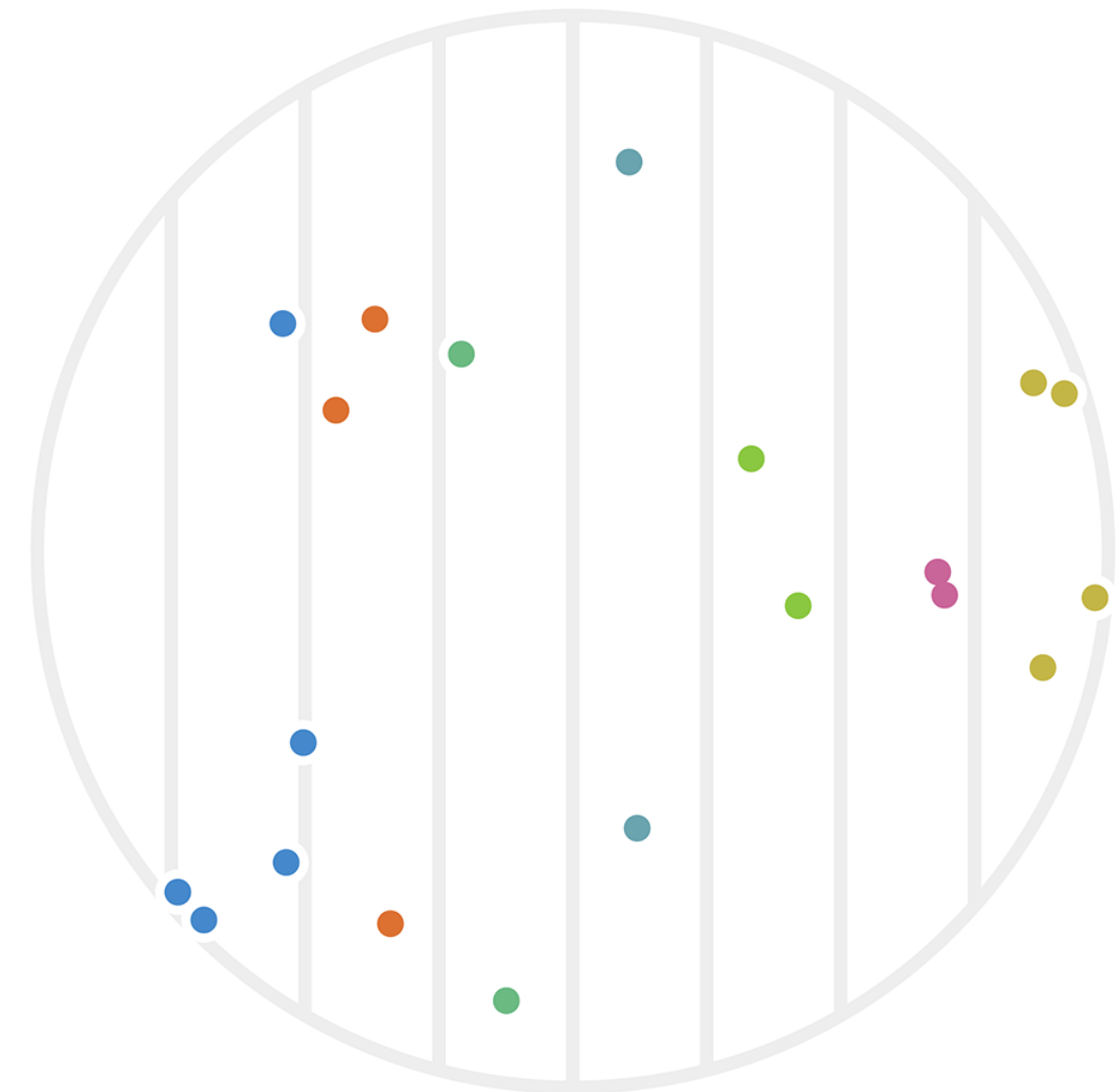- Si $d(x, y) \geq cR$ entonces $P_{\mathcal{H}}[h(x) = h(y)] \leq p_2$

# Hashing points *with projections*

**Ejemplo:**

$$h_1 : \mathbb{R}^2 \longrightarrow \mathbb{Z} \qquad \boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$$

$$h_1(\boldsymbol{x}) = \lfloor x_1 \rfloor$$
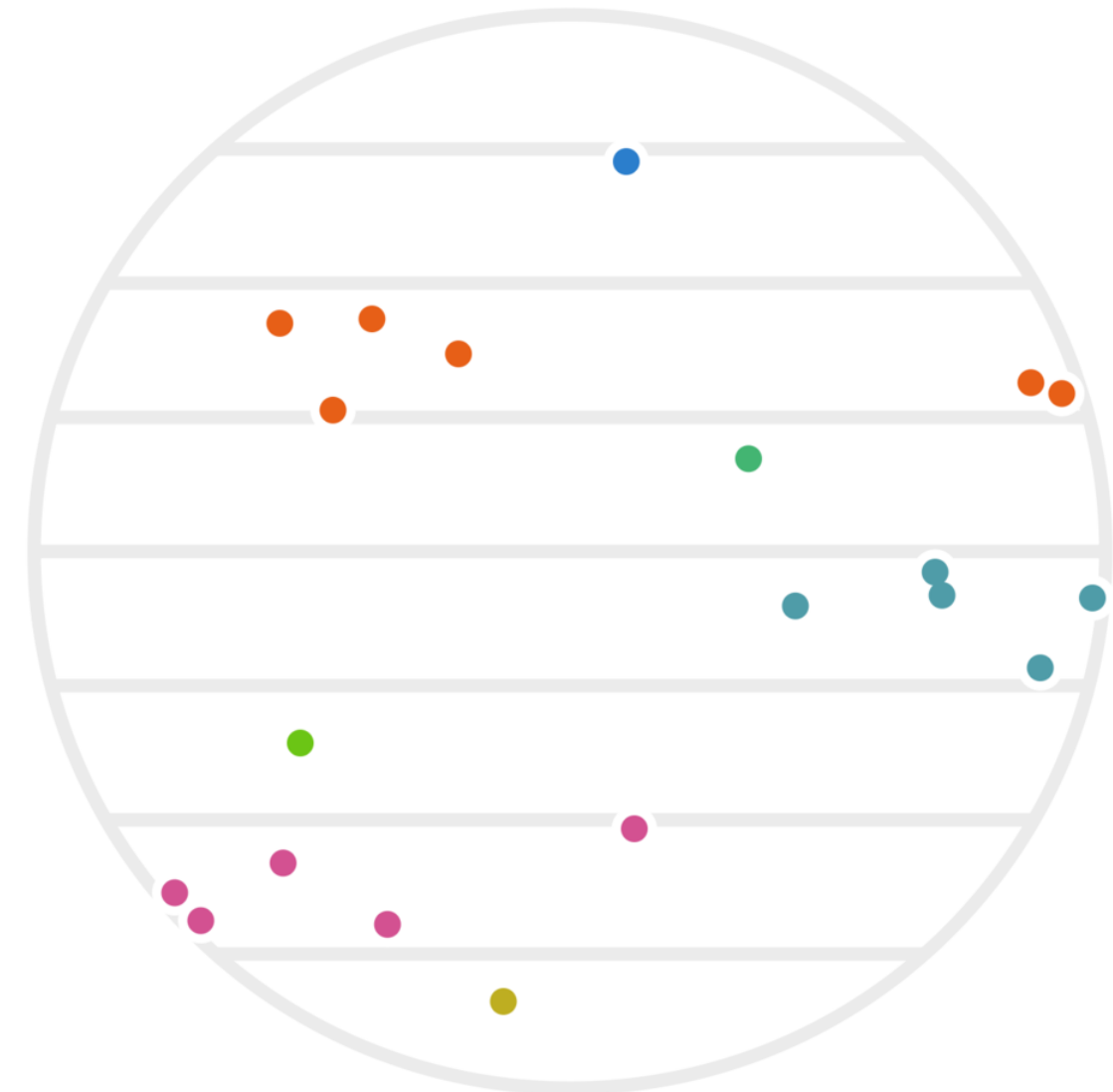
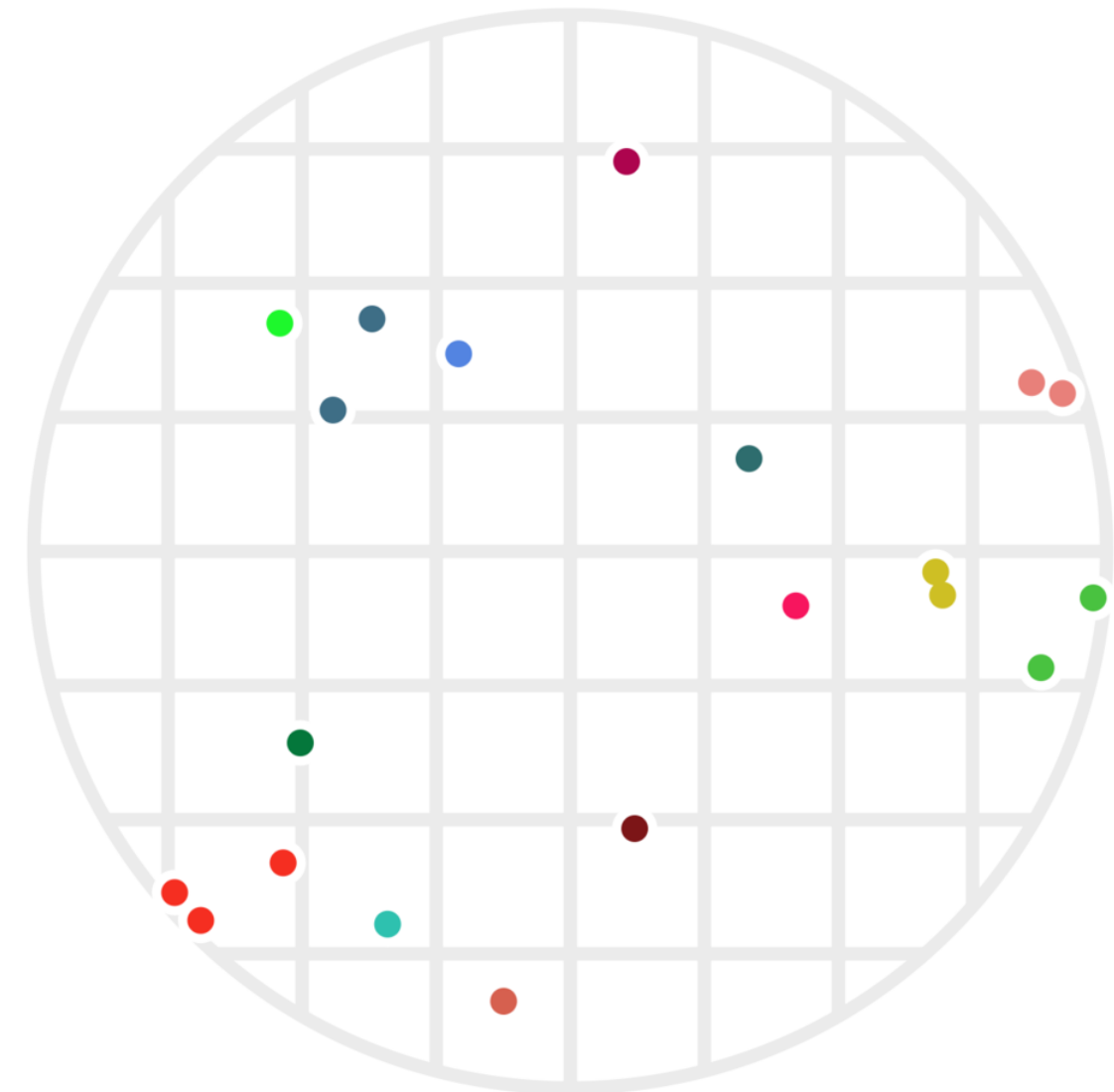

$$a \sim b \Leftrightarrow h_1(a) = h_1(b)$$

# Hashing points *with projections*

**Ejemplo:**

$h_2 : \mathbb{R}^2 \longrightarrow \mathbb{Z}$　　　　$\boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$

$h_2(\boldsymbol{x}) = \lfloor x_2 \rfloor$

$a \sim b \Longleftrightarrow h_1(a) = h_1(b)$

# Hashing points *with projections*

**Ejemplo:**

$h_1: \mathbb{R}^2 \longrightarrow \mathbb{Z}$ 　　　　$\boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$

$h_2: \mathbb{R}^2 \longrightarrow \mathbb{Z}$

$$h_1(\boldsymbol{x}) = \lfloor x_1 \rfloor$$
$$h_2(\boldsymbol{x}) = \lfloor x_2 \rfloor$$

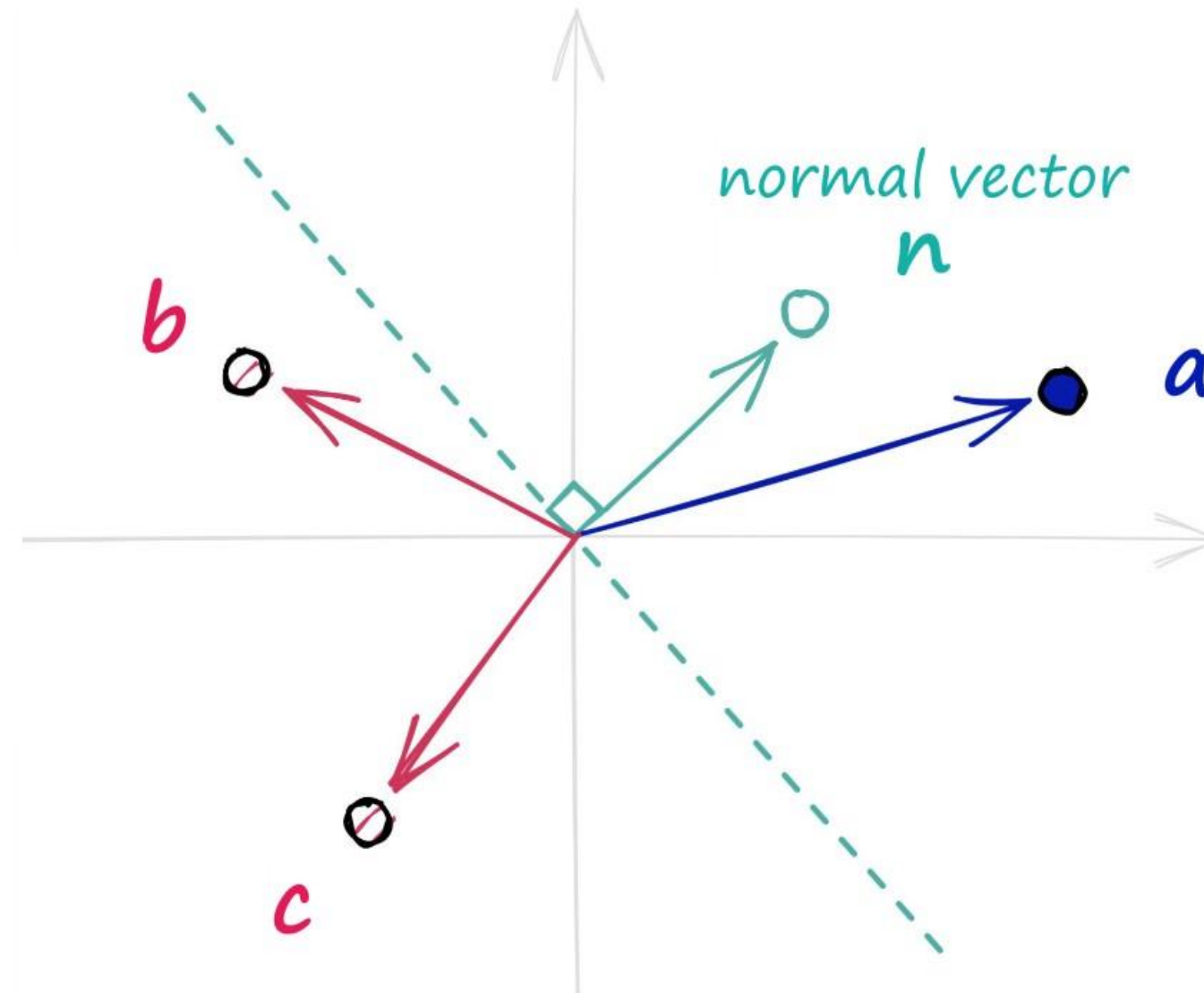$$a \sim b \Longleftrightarrow \begin{cases} h_1(a) = h_1(b) \\ h_2(a) = h_2(b) \end{cases}$$
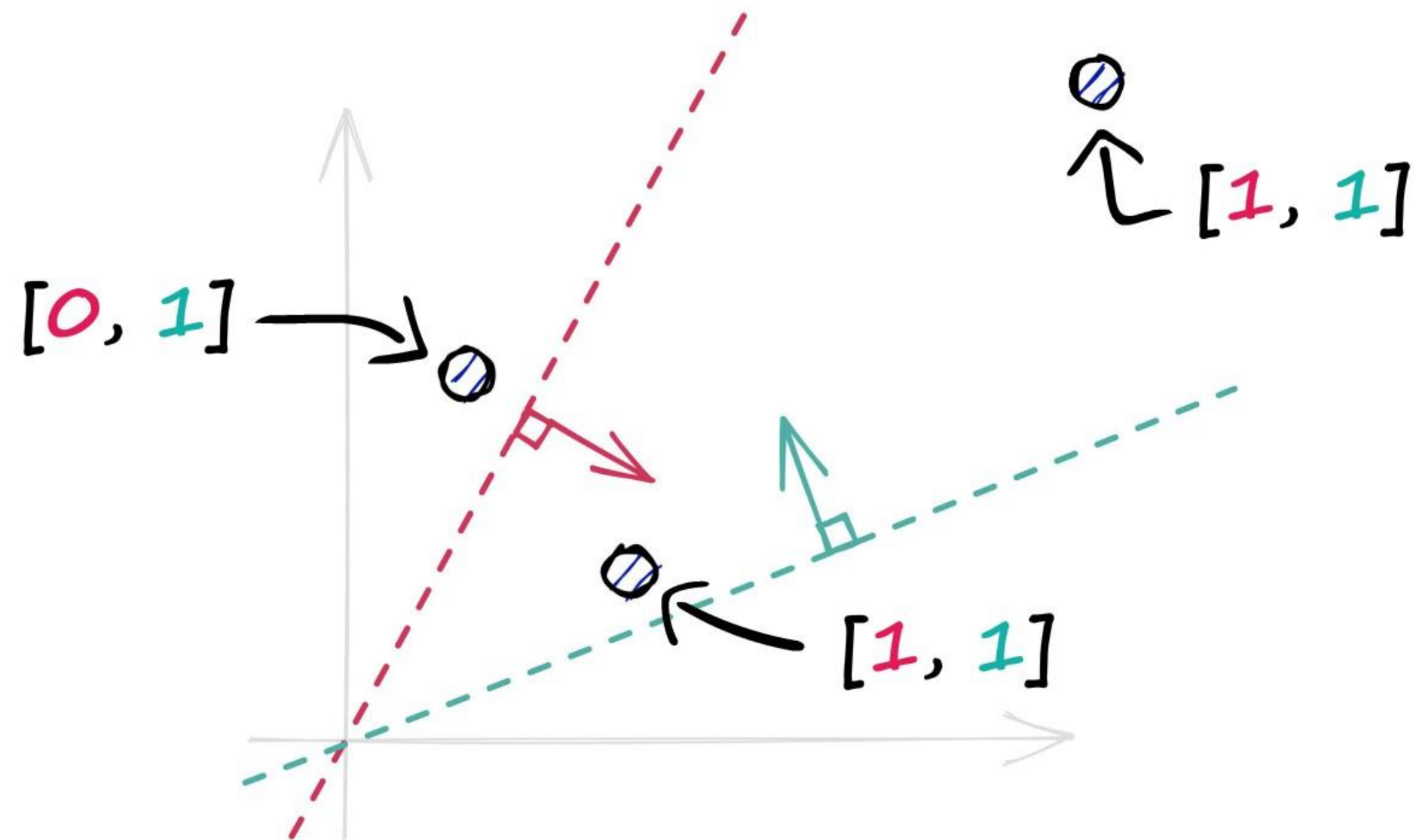
# **Hiper***planos!*

dot-product

$$n \cdot a > 0$$

$$n \cdot b < 0$$

$$n \cdot c < 0$$

normal vector

$n$

$b$

$a$

$c$

# **Hiper***planos!*

# **Proyección** *aleatoria*

**Ejemplo:**

$$h_1 \colon \mathbb{R}^2 \longrightarrow \mathbb{Z} \qquad \boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$$

$$h_1(\boldsymbol{x}) = \lfloor U x_1 + b \rfloor$$

# **Proyección** *aleatoria*

**Ejemplo:**

$$h_i: \mathbb{R}^2 \longrightarrow \mathbb{Z} \qquad \boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$$

$$h_i(\boldsymbol{x}) = \lfloor U x_i + b \rfloor$$

# **Proyección** *aleatoria*

**Ejemplo:**

$$h_i : \mathbb{R}^2 \longrightarrow \mathbb{Z} \qquad \boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$$

$$h_i(\boldsymbol{x}) = \lfloor U x_i + b \rfloor$$

# **Proyección** *aleatoria*

**Ejemplo:**

$$h_i: \mathbb{R}^2 \longrightarrow \mathbb{Z} \qquad \boldsymbol{x} = (x_1, x_2) \in \mathbb{R}^2$$

$$h_i(\boldsymbol{x}) = \lfloor U x_i + b \rfloor$$



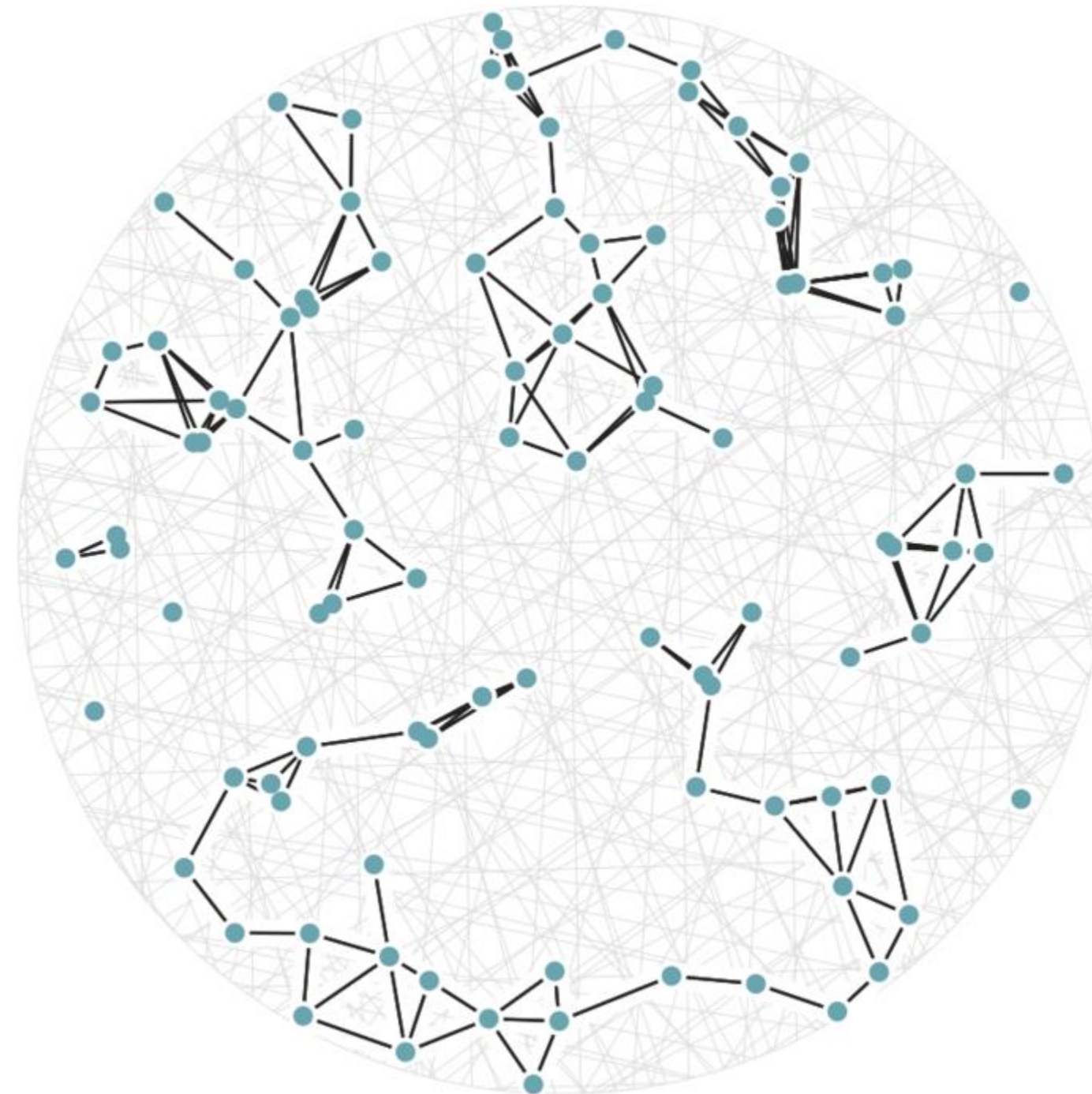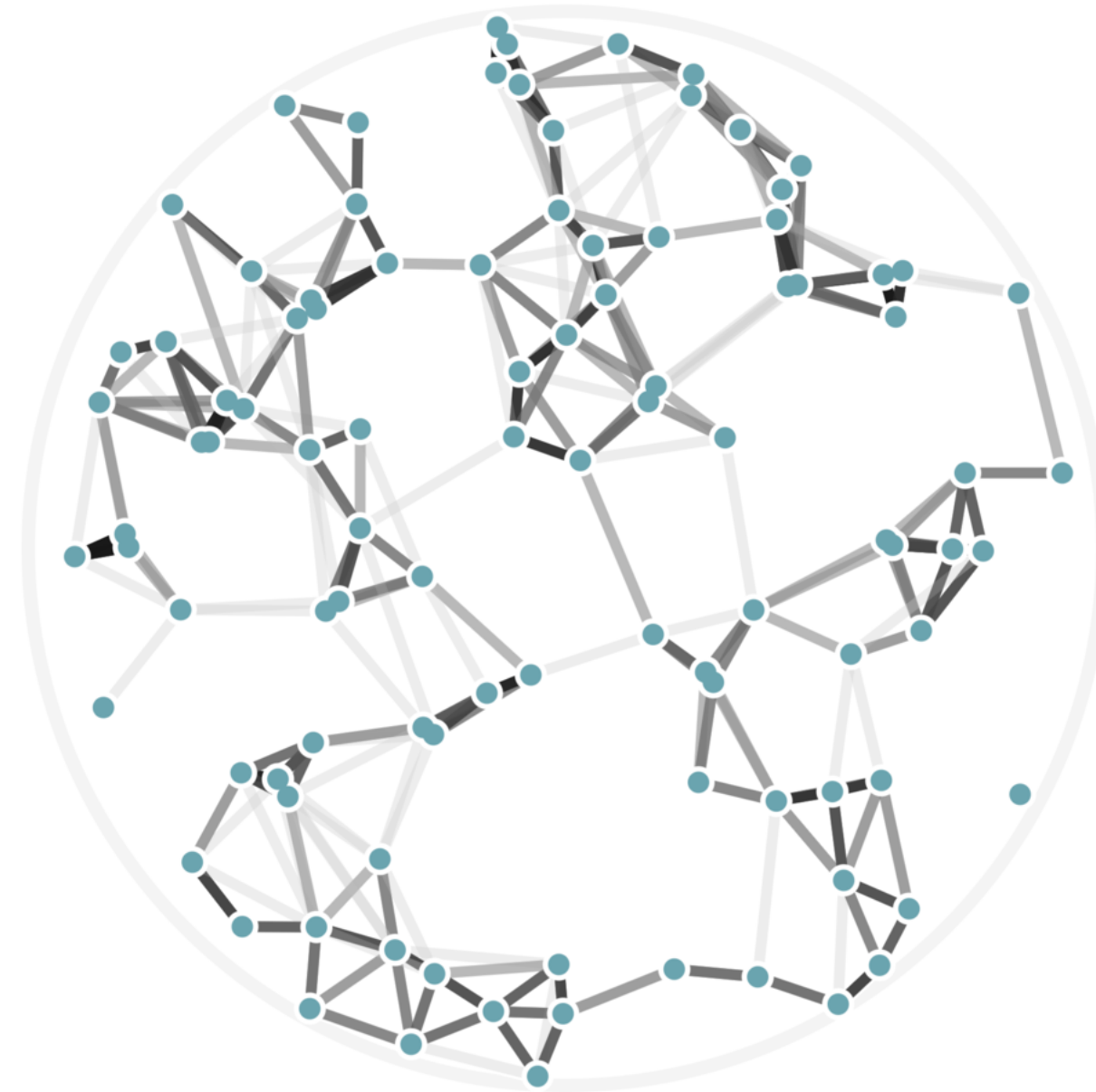$$a \sim b \Longleftrightarrow \#\{i: h_i(a) = h_i(b)\} \geq j$$

# **Proyección** *aleatoria*



$$a \sim b \Longleftrightarrow \#\{i : h_i(a) = h_i(b)\} \geq 6$$

# Proyección *aleatoria*



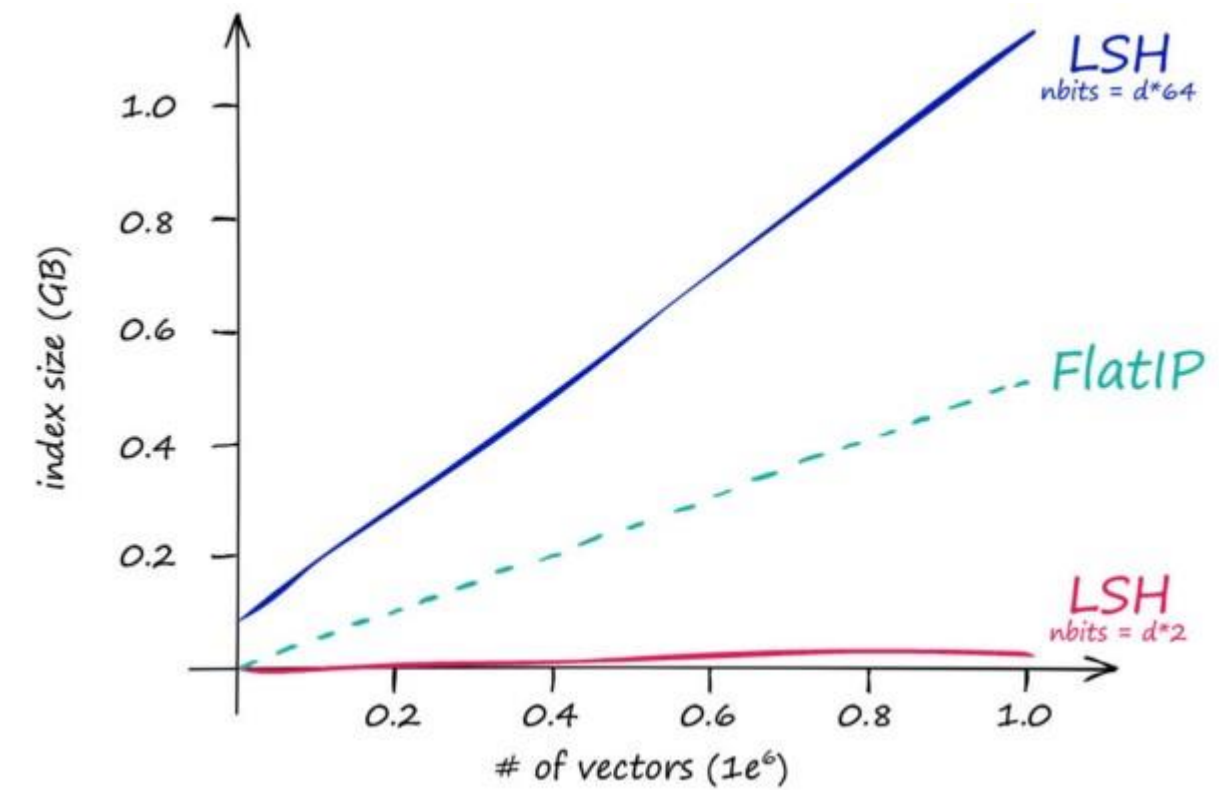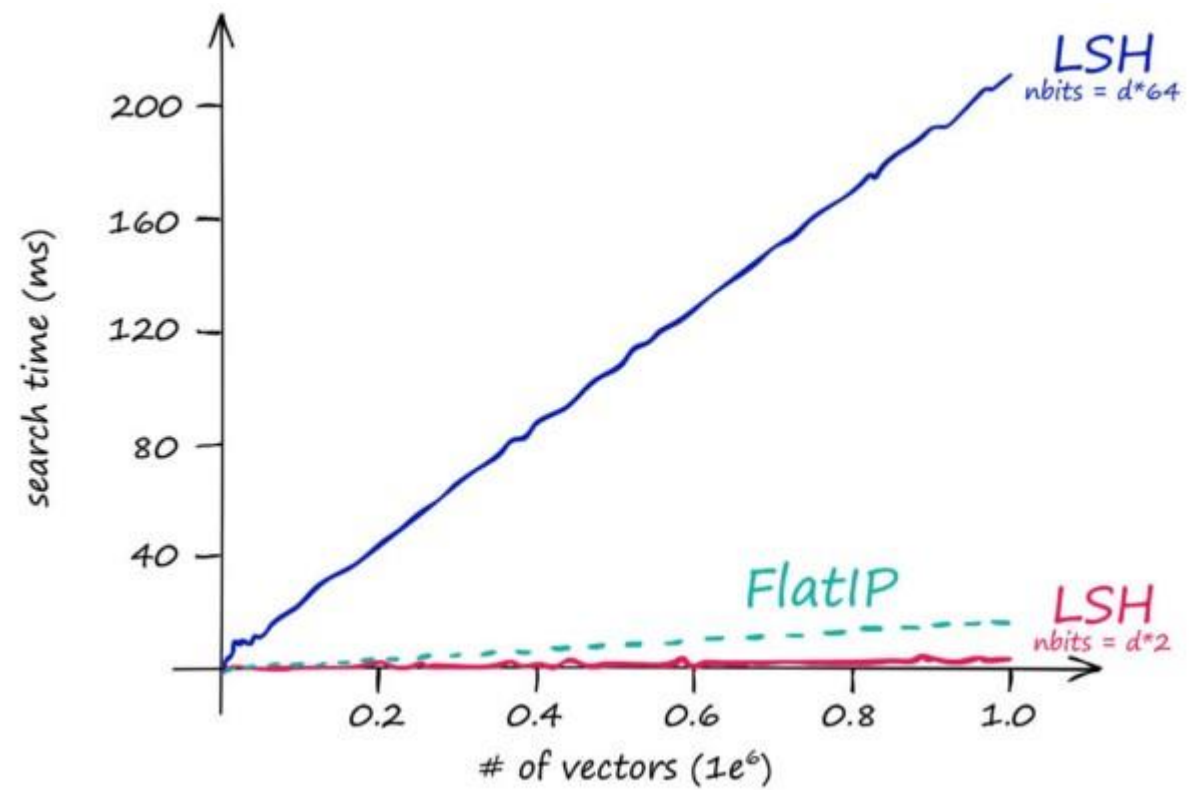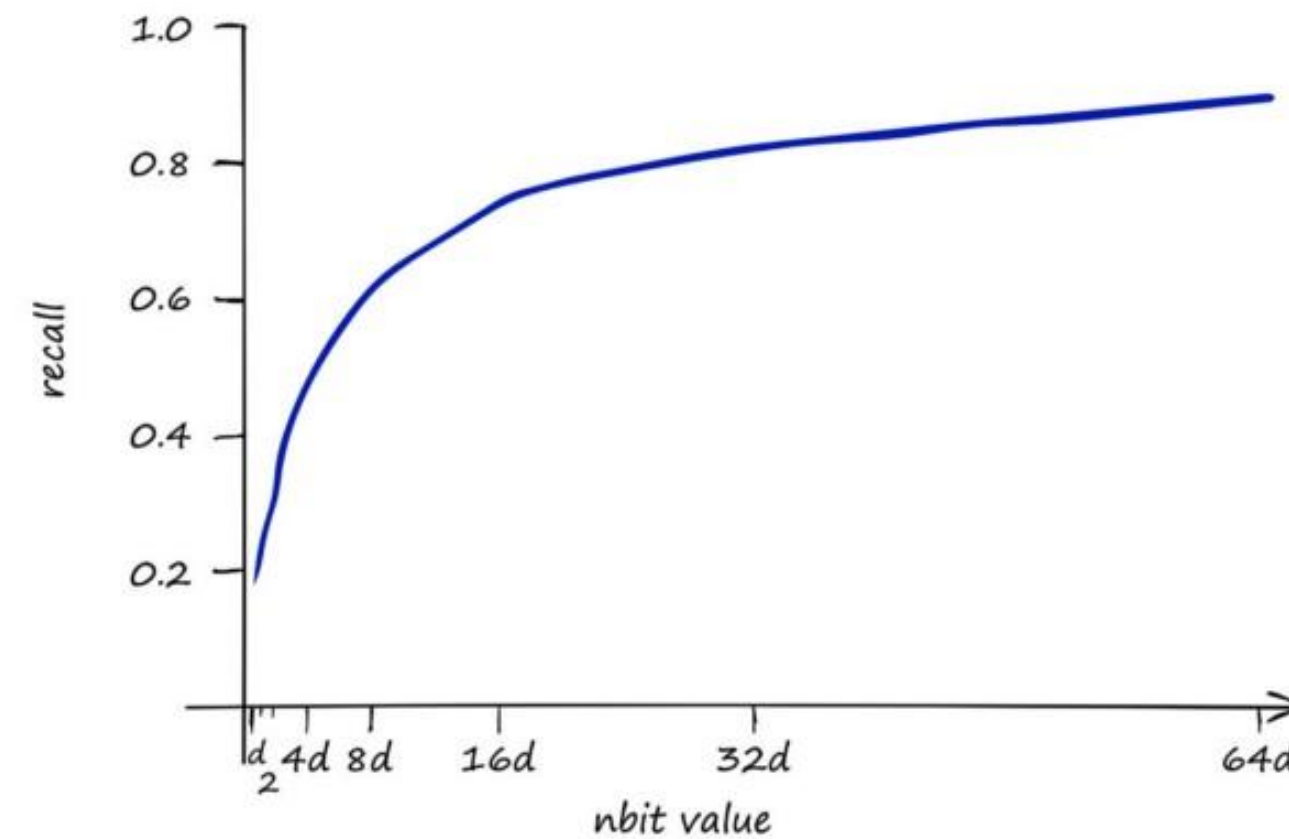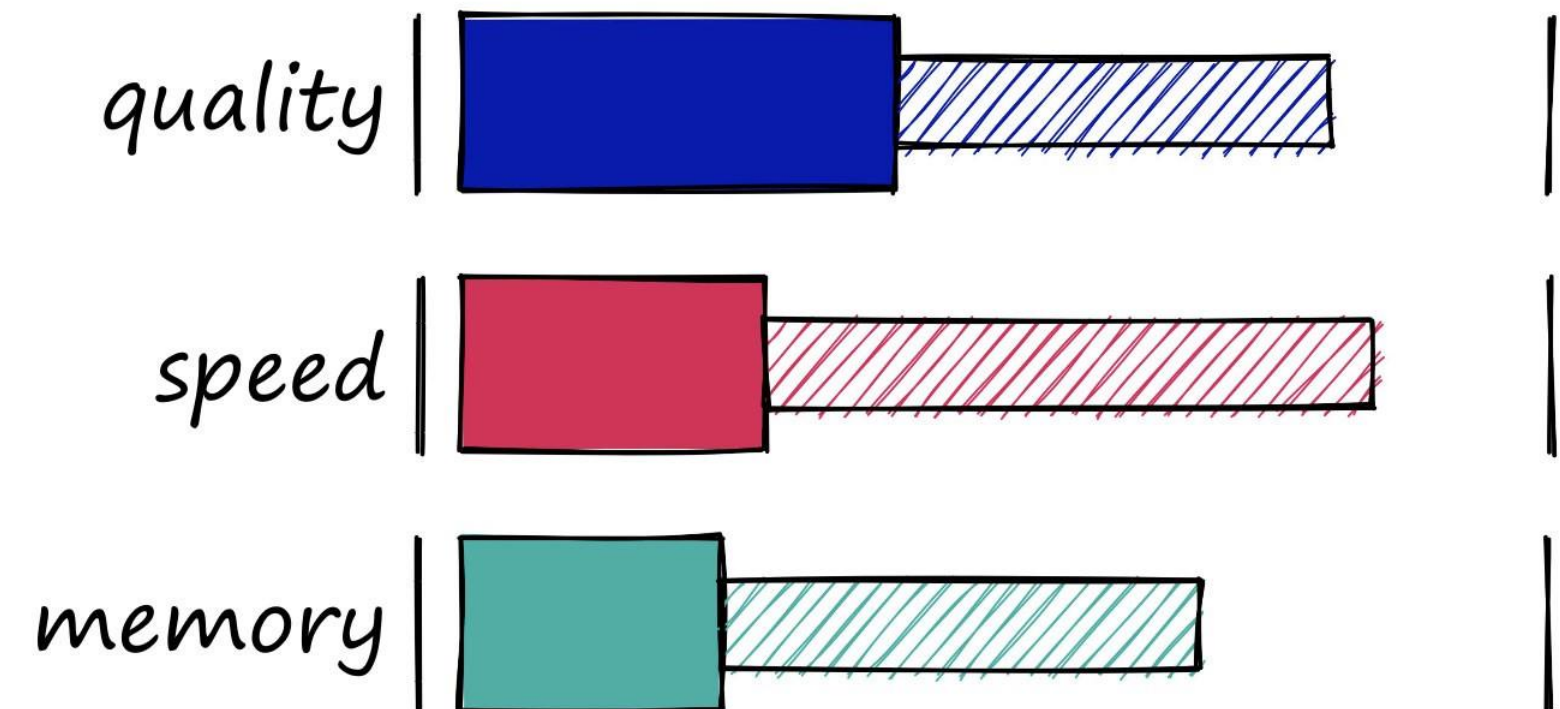$$a \sim b \iff \#\{i : h_i(a) = h_i(b)\} \geq 7$$

# Proyección *aleatoria*



$$a \sim b \iff \#\{i: h_i(a) = h_i(b)\} \geq 8$$

# **Proyección** *aleatoria*



$$a \sim b \Longleftrightarrow \#\{i : h_i(a) = h_i(b)\} \geq 9$$
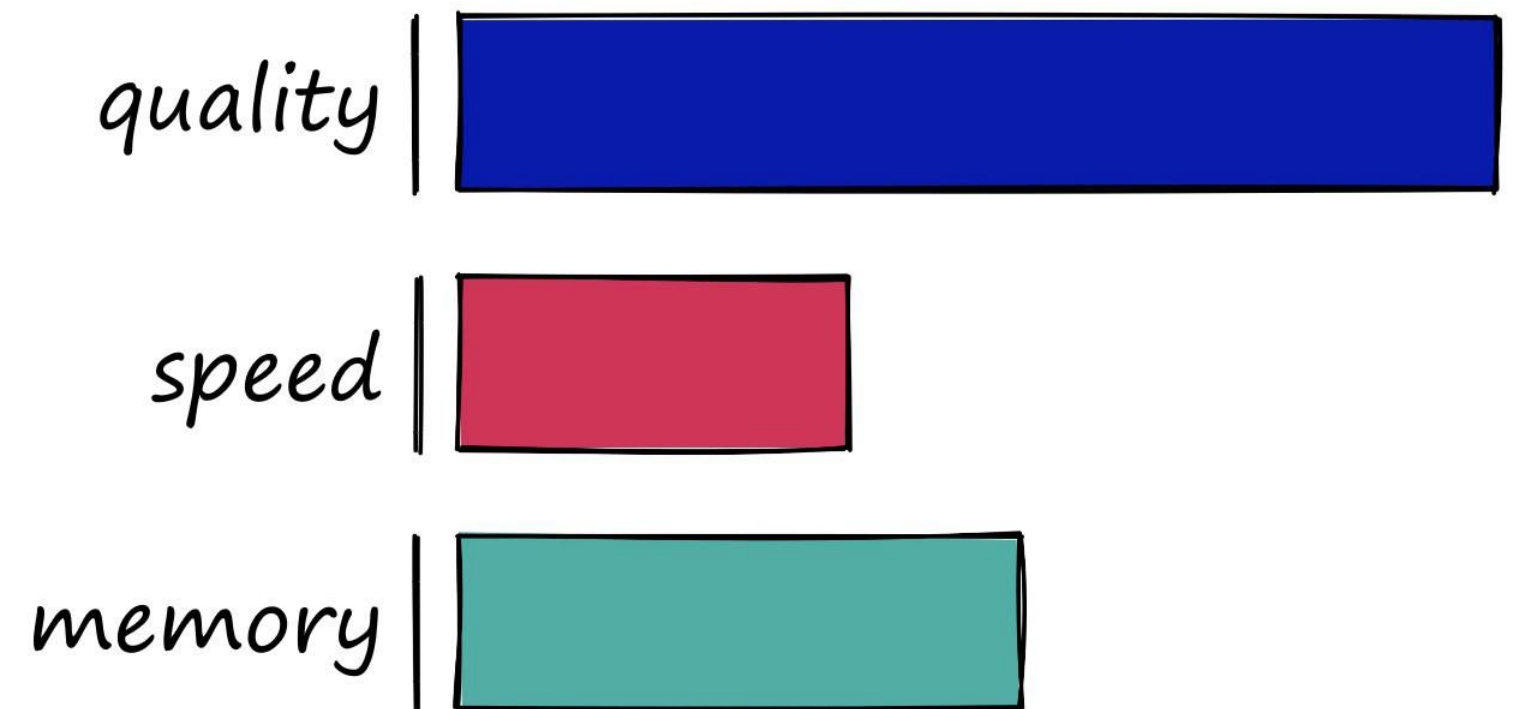
# Proyección *aleatoria*

# **Proyección** *aleatoria*
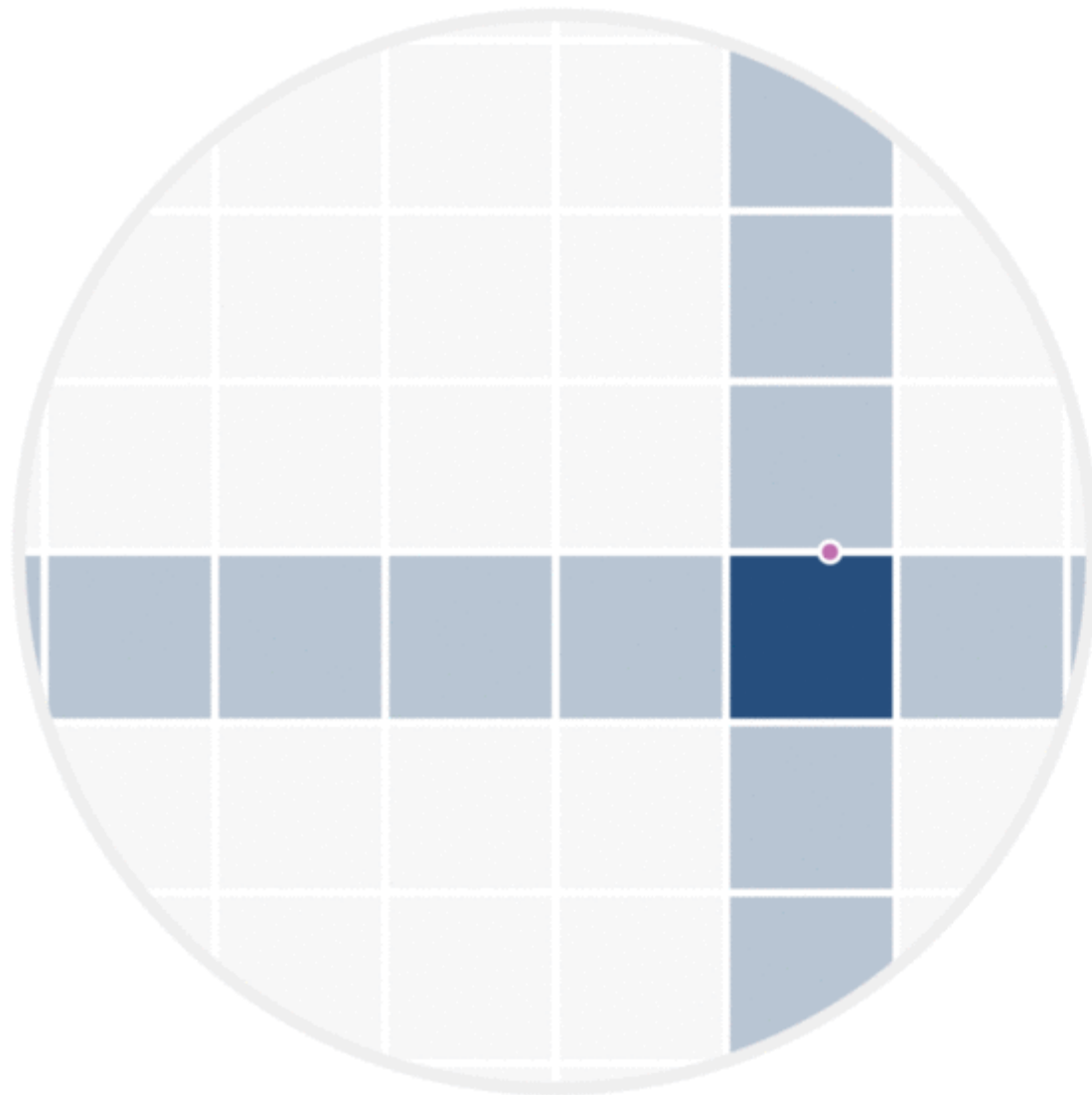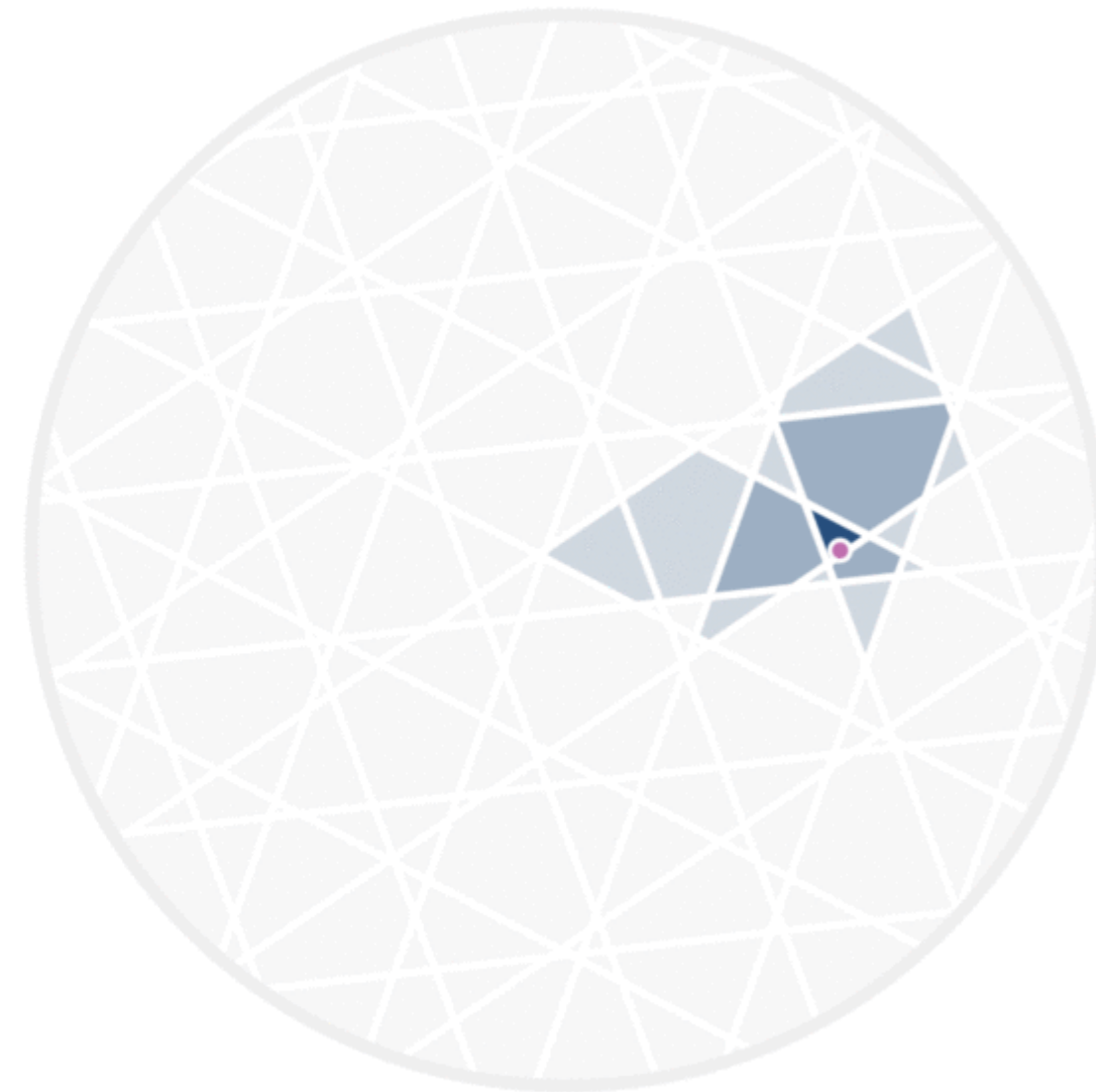
# Proyección *aleatoria*

# **Proyección** *aleatoria*
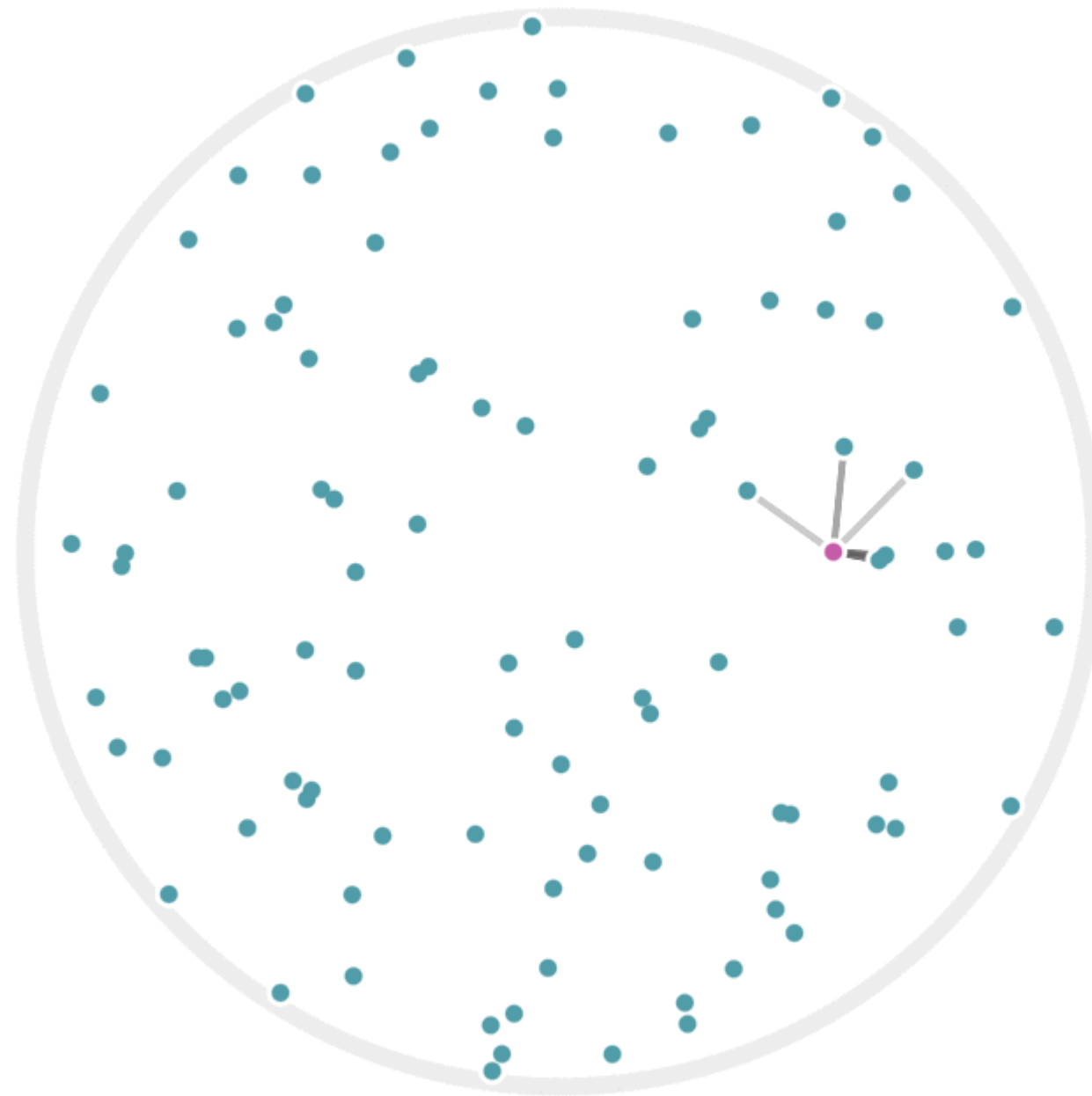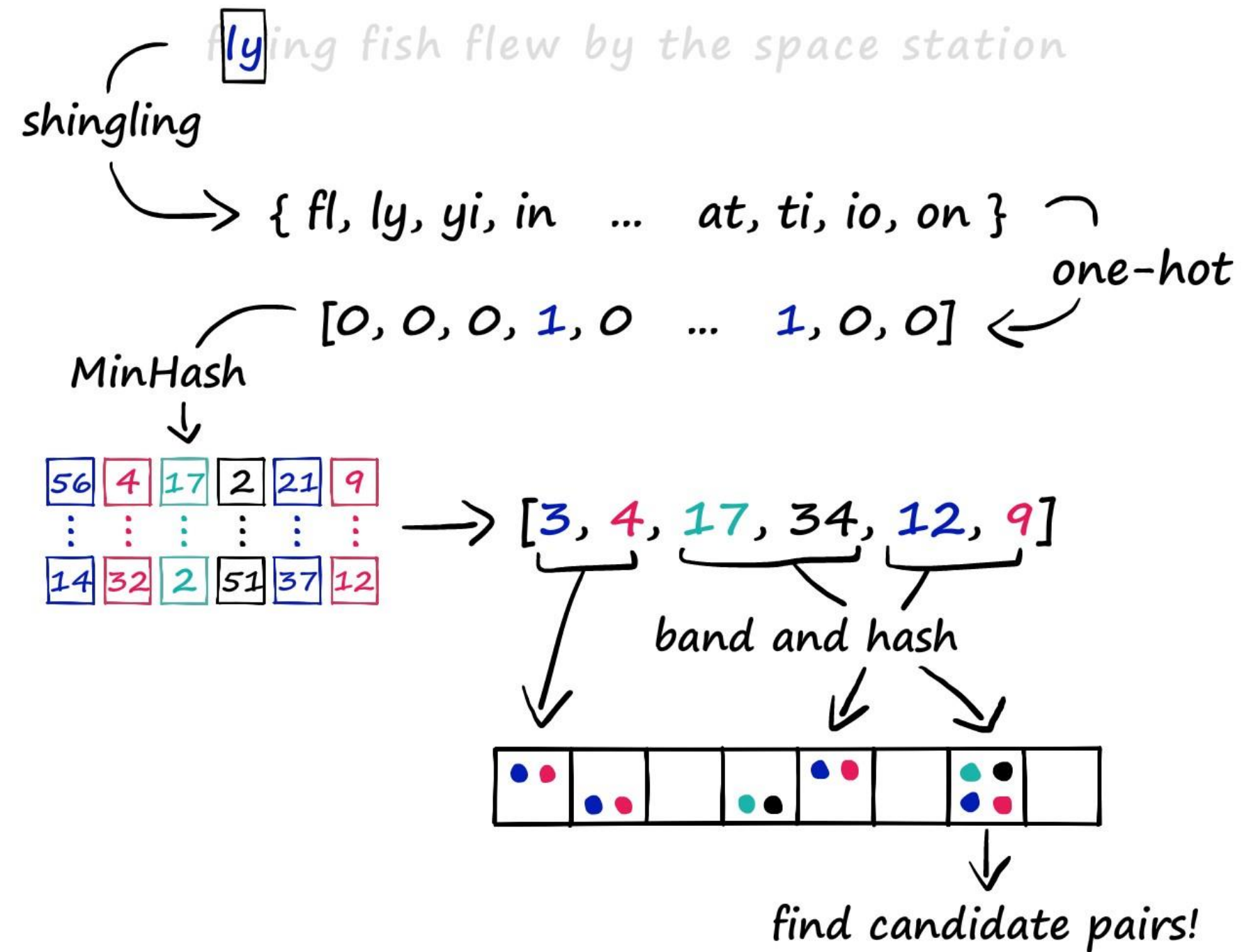
# Proyección *aleatoria*



$k = 2$

$k = 5$

# Proyección *aleatoria*

# 2.

**Min**Hashing

# Buscar frases *similares*

flying fish flew by the space station

shingling

{ fl, ly, yi, in ... at, ti, io, on }

one-hot

[0, 0, 0, 1, 0 ... 1, 0, 0]

MinHash

| 56 | 4 | 17 | 2 | 21 | 9 |
| 14 | 32 | 2 | 51 | 37 | 12 |

→ [3, 4, 17, 34, 12, 9]

band and hash

find candidate pairs!

Gabriela Hristescu and Martin Farach-Colton (1999) "Cluster-preserving embedding of proteins". Technical Report 99-50, Computer Science Department, Rutgers University

# k-*shingling*

flying fish flew by the space station

# **Vocab***ulario*

**ly**ing fish flew by the space station

...of dynamite and pet **ar**madillo

. . .

...fishing pol**e** to catch an armadillo

vocab =
{ fl, ly, yi, in    ...    di, il, ll, lo }

# **Vocab**_ulario_

# Min*Hashing*

1

0

0

1

0

1

signature:

___

# Band *method*

# **Band** *method*

# **Band** *method*

$$P = 1 - (1 - s^r)^b$$

# 3. **Sha**zam!

# **Sha**_zam_

_Avery Li-Chun Wang_ **(2003)** _"An Industrial-Strength Audio Search Algorithm"._
_En Ismir. 2003. p. 7-13._
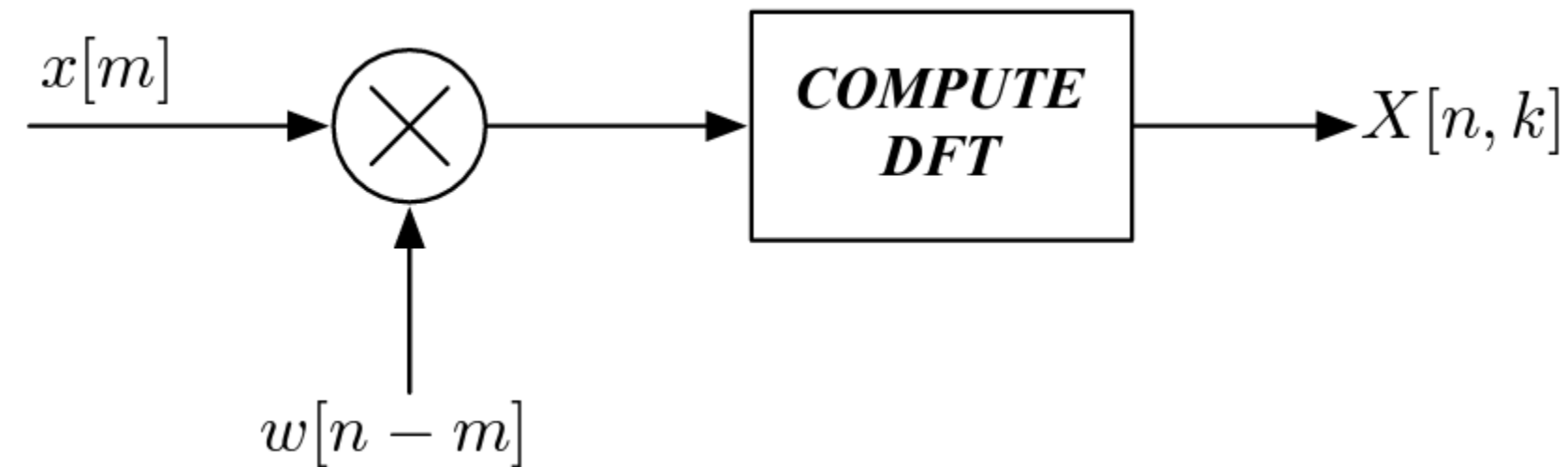
# Dominio de *Frecuencia*
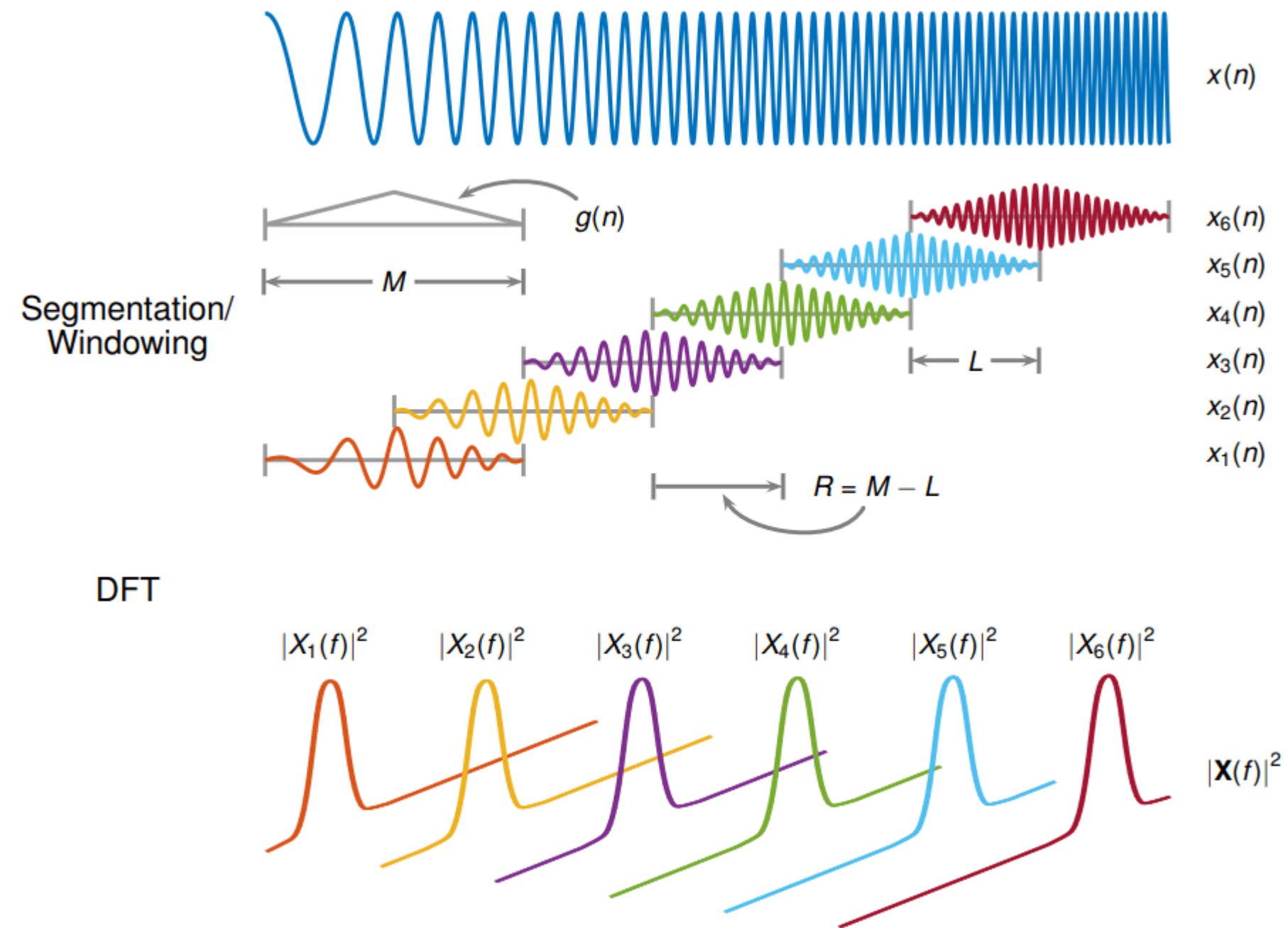
# **Short-Time** *Fourier Transform*

# **Short-Time** *Fourier Transform*


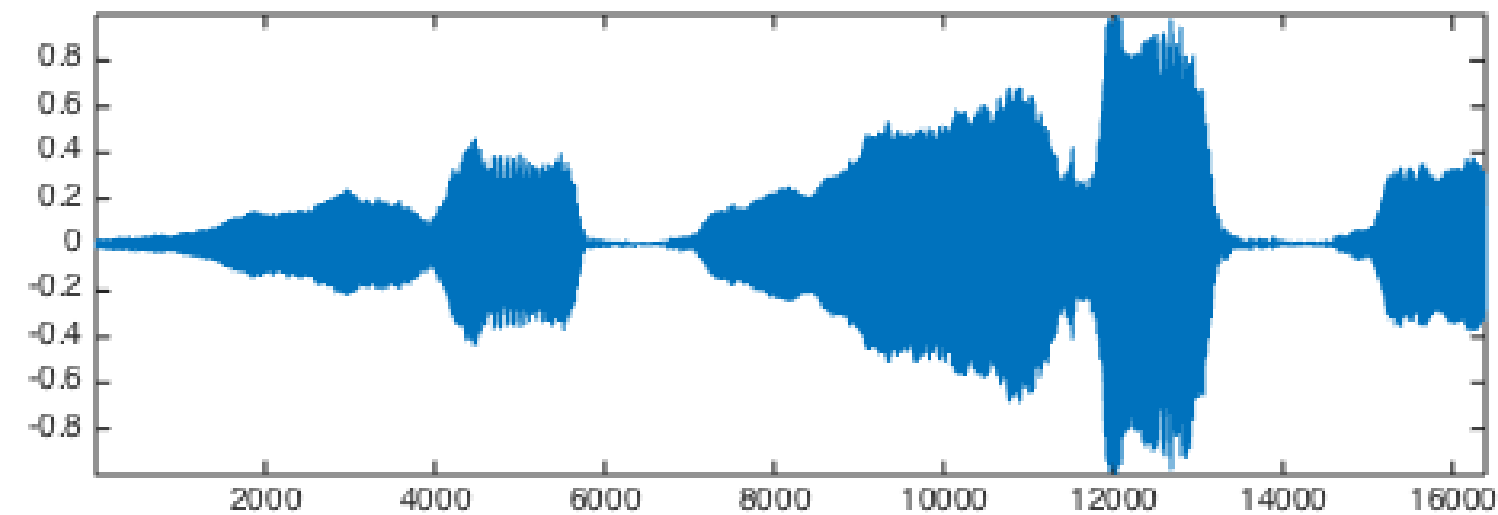
$$X[n,k] = \sum_{m=n-(N_w-1)}^{n} (x[m]w[n-m])e^{-j2\pi mk/N} = \sum_{m=n-(N_w-1)}^{n} (x[m]w[n-m])e^{-j\omega_k m}$$
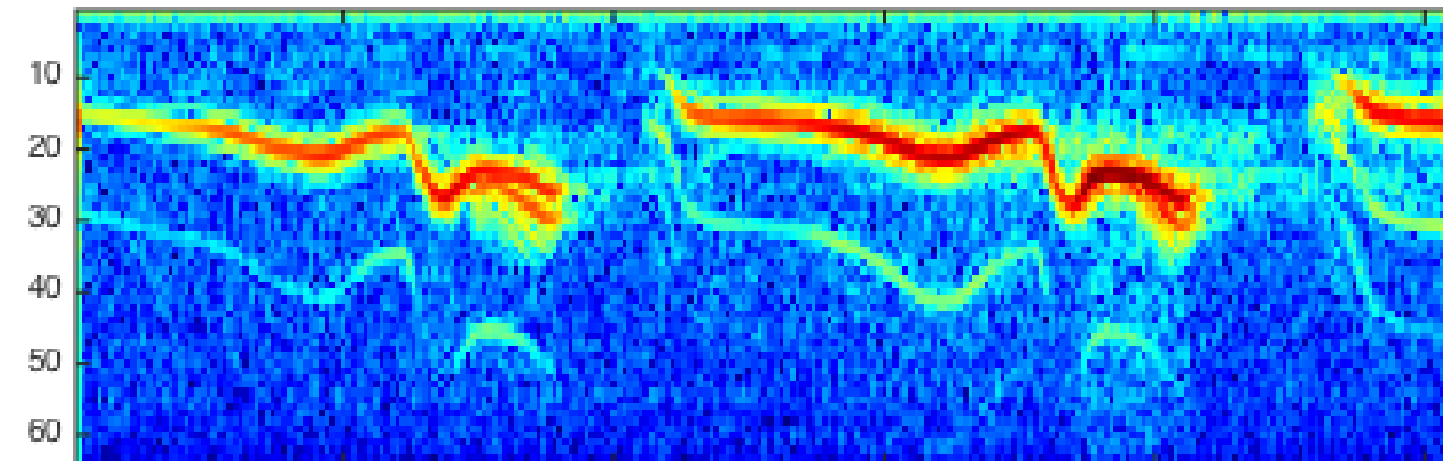
# Short-Time *Fourier Transform*

# **Short-Time** *Fourier Transform*

**Dominio de Tiempo**



**Dominio de Frecuencia**