# On Generalizing Permutation-Based Representations for Approximate Search
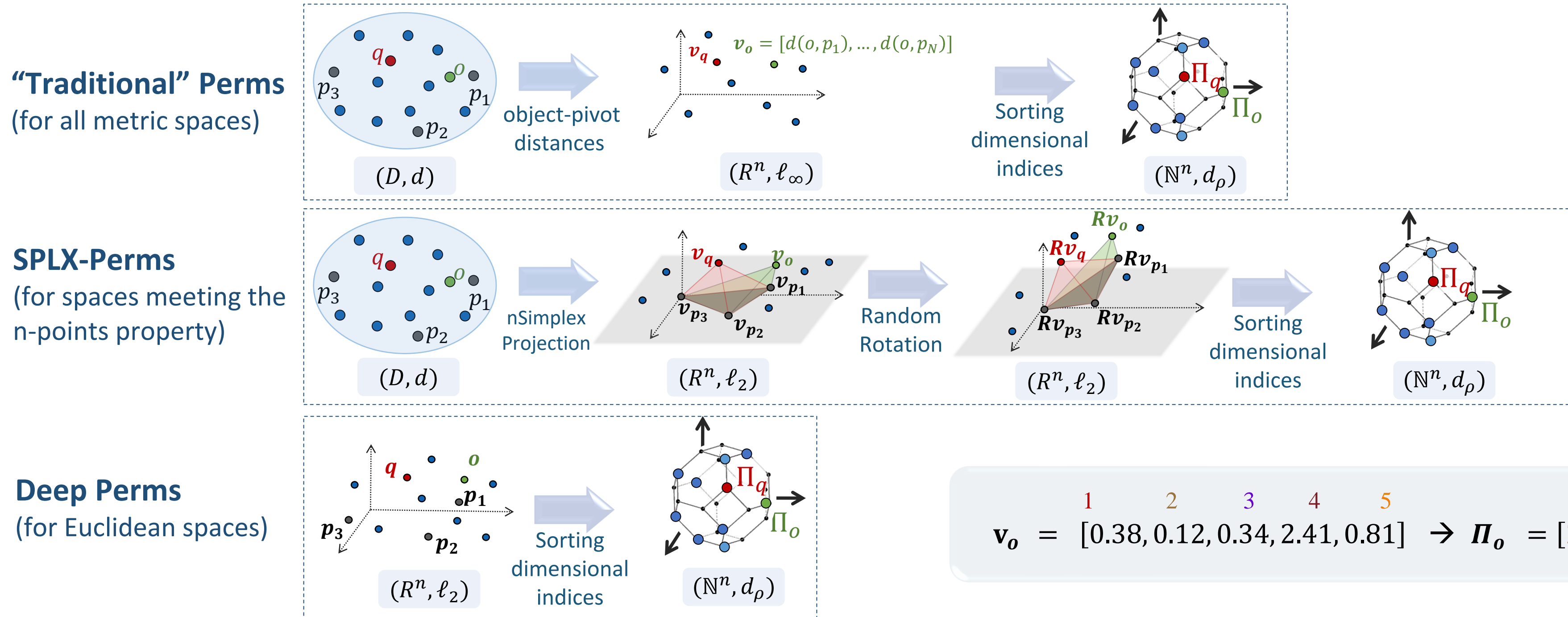
Lucia Vadicamo, Claudio Gennaro, Giuseppe Amato
*Institute of Information Science and Technologies (ISTI), CNR, Pisa, Italy*

## Background: Permutation-based methods

- Data objects are represented as **permutations** of a finite set of integers: $\Pi_o = [i_1, \ldots, i_N], \quad i_k \in \{1, \ldots, N\}$
- **Similarity queries are executed in the permutation space**
- Permutations can be **efficiently indexed and searched** (e.g., using inverted files)

**"Traditional" Perms**
(for all metric spaces)

$v_o = [d(o,p_1), \ldots, d(o,p_N)]$

object-pivot distances → $(R^n, \ell_\infty)$ → Sorting dimensional indices → $(\mathbb{N}^n, d_\rho)$

$(D,d)$

**SPLX-Perms**
(for spaces meeting the n-points property)

$(D,d)$ → nSimplex Projection → $(R^n, \ell_2)$ → Random Rotation → $(R^n, \ell_2)$ → Sorting dimensional indices → $(\mathbb{N}^n, d_\rho)$

**Deep Perms**
(for Euclidean spaces)

$(R^n, \ell_2)$ → Sorting dimensional indices → $(\mathbb{N}^n, d_\rho)$

$$v_o = [0.38, 0.12, 0.34, 2.41, 0.81] \rightarrow \Pi_o = [2,3,1,5,4]$$

## Our generalization: Permutations induced by a space transformation $f$

**Def.** The *permutation representation* of an object $o \in (D,d)$ induced by the function $f : (D,d) \to \mathbb{R}^N$ is the sequence $\Pi_o^f = [\pi_1, \ldots, \pi_N]$ that lists the *permutants* $\{1, \ldots, N\}$ in an order such that $\forall\, i \in \{1, \cdots, N-1\}$
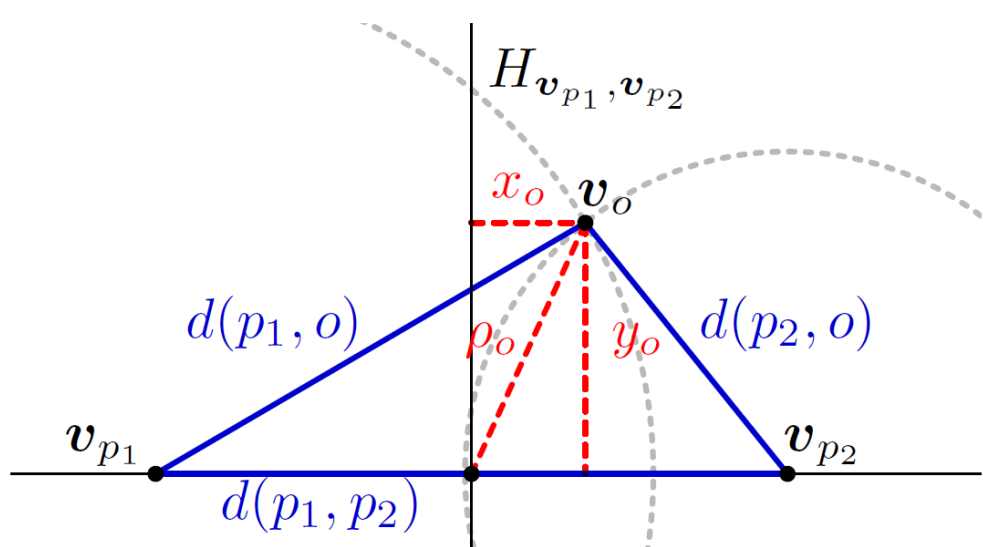
$$f(o)_{\pi_i} < f(o)_{\pi_{i+1}} \quad \text{or} \quad \left[ f(o)_{\pi_i} = f(o)_{\pi_{i+1}} \right] \wedge [\pi_i < \pi_{i+1}]$$

Permutants: 1 2 3 4 5
$$o \in (D,d) \to f(o) = [0.3, 0.1, 0.4, 2.4, 1.1] \in \mathbb{R}^N \to sort(f(o)) = [0.1, 0.3, 0.4, 1.1, 2.4] \to \Pi_o = [2,1,3,5,4]$$

➢ *novel permutation-based representations can be defined (assuming a suitable $f : (D,d) \to \mathbb{R}^N$ is used!)*

## Pivot Pair Permutations

**Basic idea using 2 pivots:**

$$\rho_o = \sqrt{x_o^2 + y_o^2} = \frac{1}{2}\sqrt{2d(o,p_1) + 2d(o,p_2) - d(p_1,p_2)^2}$$

$\boldsymbol{\rho_o}$ can be interpreted as the **distance to an artificial pivot** that is equidistant to the two original pivots

**Using $n$ pivots $\{p_1, \ldots, p_n\}$ and $m$ pivot pairs $(p_{i_1}, p_{i_2})$, $i = 1, \ldots, m$:**

$$f' : (D,d) \to \mathbb{R}^m$$
$$o \to \left[\rho_o^{(1)}, \ldots, \rho_o^{(m)}\right]$$

Sorting dimensional indices → $\Pi_o^{f'}$ → **Pairs Permutation** (P-Perm)

$$f'' : (D,d) \to \mathbb{R}^{n+m}$$
$$o \to \left[d(o,p_1), \ldots, d(o,p_n), \rho_o^{(1)}, \ldots, \rho_o^{(m)}\right]$$

$\Pi_o^{f''}$ → **Pivot-Pairs Permutation** (PP-Perm)

**PP-Perms:** best trade-off between recall, search cost and the cost for computing the permutations

**Future work:** what properties should a function $f : (D,d) \to \mathbb{R}^N$ satisfy to produce good permutations for approximate search? Can we learn $f$?

**Clustered Euclid30**
- Perms (n=5000)
- PP-Perms (n=500 + m=4500)
- P-Perms (m=5000), n=500
- Perms (n=500)
- Sequential Scan

Recall@10 vs Search Cost (MB), (log scale)

**Gaussian Euclid30**
- Perms (n=5000)
- PP-Perms (n=500 + m=4500)
- P-Perms (m=5000), n=500
- Perms (n=500)
- Sequential Scan

Recall@10 vs Search Cost (MB), (log scale)

**ANN-SIFT**
- Perms (n=10K)
- PP-Perms (n=1K + m=9K)
- Perms (n=1K)
- Perms (n=10K), re-rank
- PP-Perms (n=1K + m=9K), re-rank
- Perms (n=1K), re-rank
- Sequential Scan (original dataset)

Recall@10 vs Search Cost (MB), (log scale)

**CoPhIR**
- Perms (n=10K)
- PP-Perms (n=1K + m=9K)
- Perms (n=1K)
- Perms (n=10K), re-rank
- PP-Perms (n=1K + m=9K), re-rank
- Perms (n=1K), re-rank
- Sequential Scan (original dataset)

Recall@10 vs Search Cost (MB), (log scale)