

Approximate Conditional Coverage via Neural Model Approximations

Allen Schmaltz and Danielle Rasooly

Reexpress AI, Inc. and Harvard University

Workshop on Distribution-Free Uncertainty Quantification at ICML (July 2022), Baltimore, Maryland


5 minute overview

Desiderata for deploying deep learning models

- High accuracy for point predictions
 - Typically necessary but *not sufficient*

Desiderata for deploying deep learning models

- High accuracy for point predictions
 - Typically necessary but *not sufficient*
- Interpretability
 - Map: Relevant feature subsets to those with known labels
- Local updatability
- Uncertainty quantification



We seek a general framework with all of these properties

Desiderata for uncertainty quantification

- Minimize distributional assumptions
- Robustness to data shifts (covariate & label shift)
- Out-of-distribution data does not lead to unexpected catastrophes
- Minimize free parameters (i.e., practical reliability)
- Informative coverage
 - Prefer conservative coverage over under-coverage

Split-conformal prediction sets for classification

- Computationally expensive blackbox: F
- Training dataset: $\mathcal{D}_{\text{tr}} = \{(X_i, Y_i)\}_{i=1}^I$ with $Y_i \in \mathcal{Y} = \{1, \dots, C\}$
- Held-out labeled calibration dataset: $\mathcal{D}_{\text{ca}} = \{(X_j, Y_j)\}_{j=I+1}^{N=I+J}$

Split-conformal prediction sets for classification

- Computationally expensive blackbox: F
- Training dataset: $\mathcal{D}_{\text{tr}} = \{(X_i, Y_i)\}_{i=1}^I$ with $Y_i \in \mathcal{Y} = \{1, \dots, C\}$
- Held-out labeled calibration dataset: $\mathcal{D}_{\text{ca}} = \{(X_j, Y_j)\}_{j=I+1}^{N=I+J}$
- **Seek:** A prediction set $\hat{\mathcal{C}}(X_{N+1}) \in 2^{\mathcal{C}}$ for a new, unseen test instance X_{N+1} from \mathcal{D}_{te}
 - Contains the true label with coverage level $1 - \alpha \in (0, 1)$ *on average*

Split-conformal prediction sets for classification

- Computationally expensive blackbox: F
- Training dataset: $\mathcal{D}_{\text{tr}} = \{(X_i, Y_i)\}_{i=1}^I$ with $Y_i \in \mathcal{Y} = \{1, \dots, C\}$
- Held-out labeled calibration dataset: $\mathcal{D}_{\text{ca}} = \{(X_j, Y_j)\}_{j=I+1}^{N=I+J}$
- **Seek:** A prediction set $\hat{\mathcal{C}}(X_{N+1}) \in 2^{\mathcal{C}}$ for a new, unseen test instance X_{N+1} from \mathcal{D}_{te}
 - Contains the true label with coverage level $1 - \alpha \in (0, 1)$ *on average*
- Finite-sample *marginal* guarantee:
 - $\mathbb{P} \left\{ Y_{N+1} \in \hat{\mathcal{C}}(X_{N+1}) \right\} \geq 1 - \alpha$
 - Via $\hat{\mathcal{C}}(x_{N+1}) = \{c \in \mathcal{Y} : \hat{\pi}^c(x_{N+1}) \geq \hat{\tau}^\alpha\}$, where $\hat{\tau}^\alpha = 1 - \hat{l}^\alpha$

Quantile threshold

conditional coverage

- Finite-sample *conditional* coverage:

$$\mathbb{P} \left\{ Y_{N+1} \in \hat{\mathcal{C}}(X_{N+1}) \mid X_{N+1} = x \right\} \geq 1 - \alpha$$

conditional coverage

- Finite-sample *conditional* coverage:

$$\mathbb{P} \left\{ Y_{N+1} \in \hat{\mathcal{C}}(X_{N+1}) \mid X_{N+1} = x \right\} \geq 1 - \alpha$$



Approximate conditional coverage

- Finite-sample *conditional* coverage:

$$\mathbb{P} \left\{ Y_{N+1} \in \hat{\mathcal{C}}(X_{N+1}) \mid X_{N+1} = x \right\} \geq 1 - \alpha$$

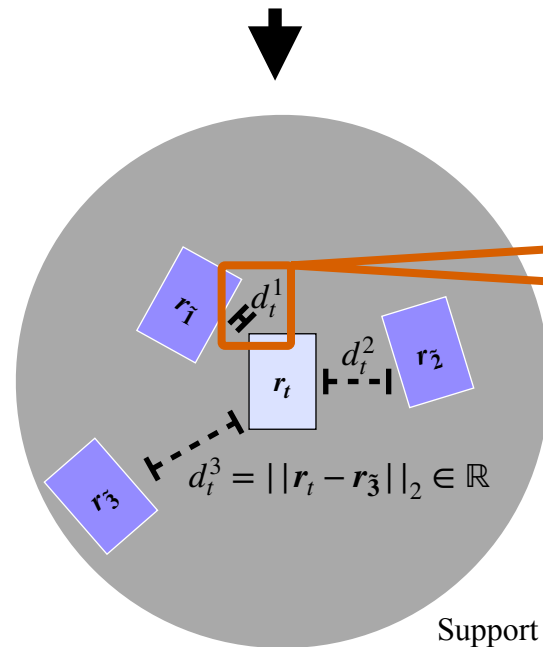
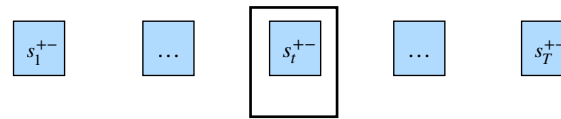


- Finite-sample *approximate conditional* coverage:

$$\mathbb{P} \left\{ Y_{N+1} \in \hat{\mathcal{C}}(X_{N+1}) \mid X_{N+1} \in \mathcal{B}(x), Y_{N+1} = y \right\} \geq 1 - \alpha, \text{ with } P_X(\mathcal{B}(x)) \geq \xi$$

Recast a prediction as a weighting over the training set

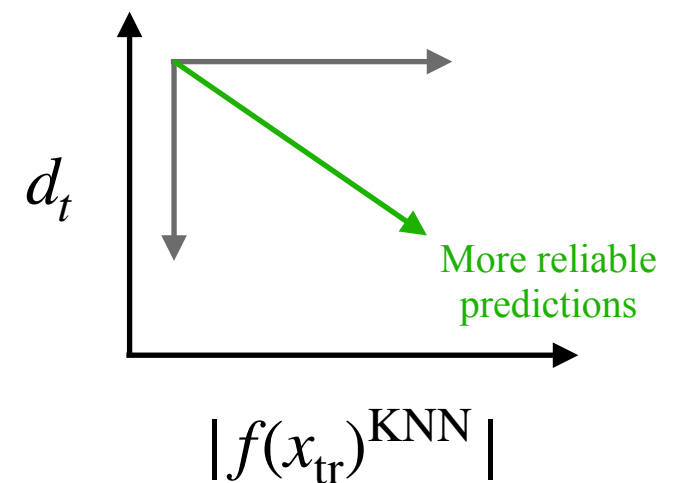
KNN Approximation



Model uncertainty: This bounded value reaches its min/max when $\tanh(s_k^{+-})$ & y (or y_k , with token-level labels) agree, for all k (assuming $\gamma > 0$).

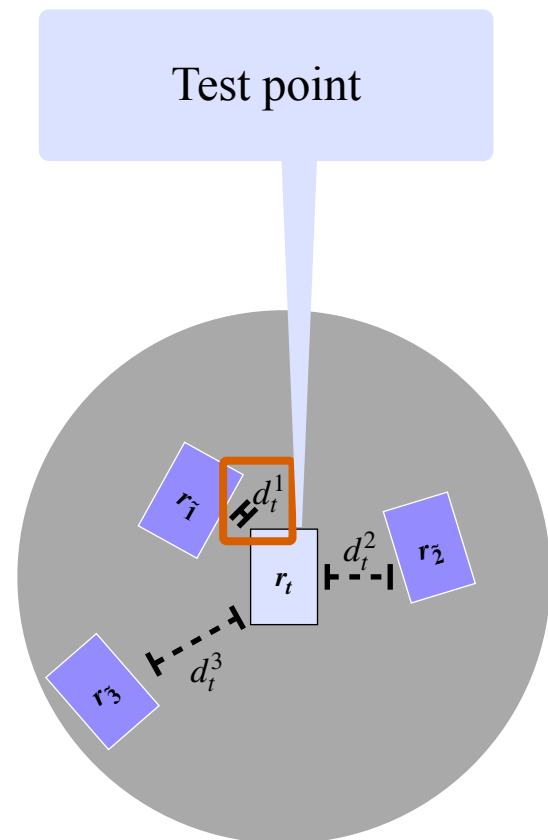
$$s_t^{+-} \approx \beta + w_1 \cdot (\tanh(s_1^{+-}) + \gamma \cdot y_1) + w_2 \cdot (\tanh(s_2^{+-}) + \gamma \cdot y_2) + w_3 \cdot (\tanh(s_3^{+-}) + \gamma \cdot y_3)$$

$$w_k = \frac{\exp(-d_k/\eta)}{\sum_{k'=1}^3 \exp(-d_{k'}/\eta)}$$

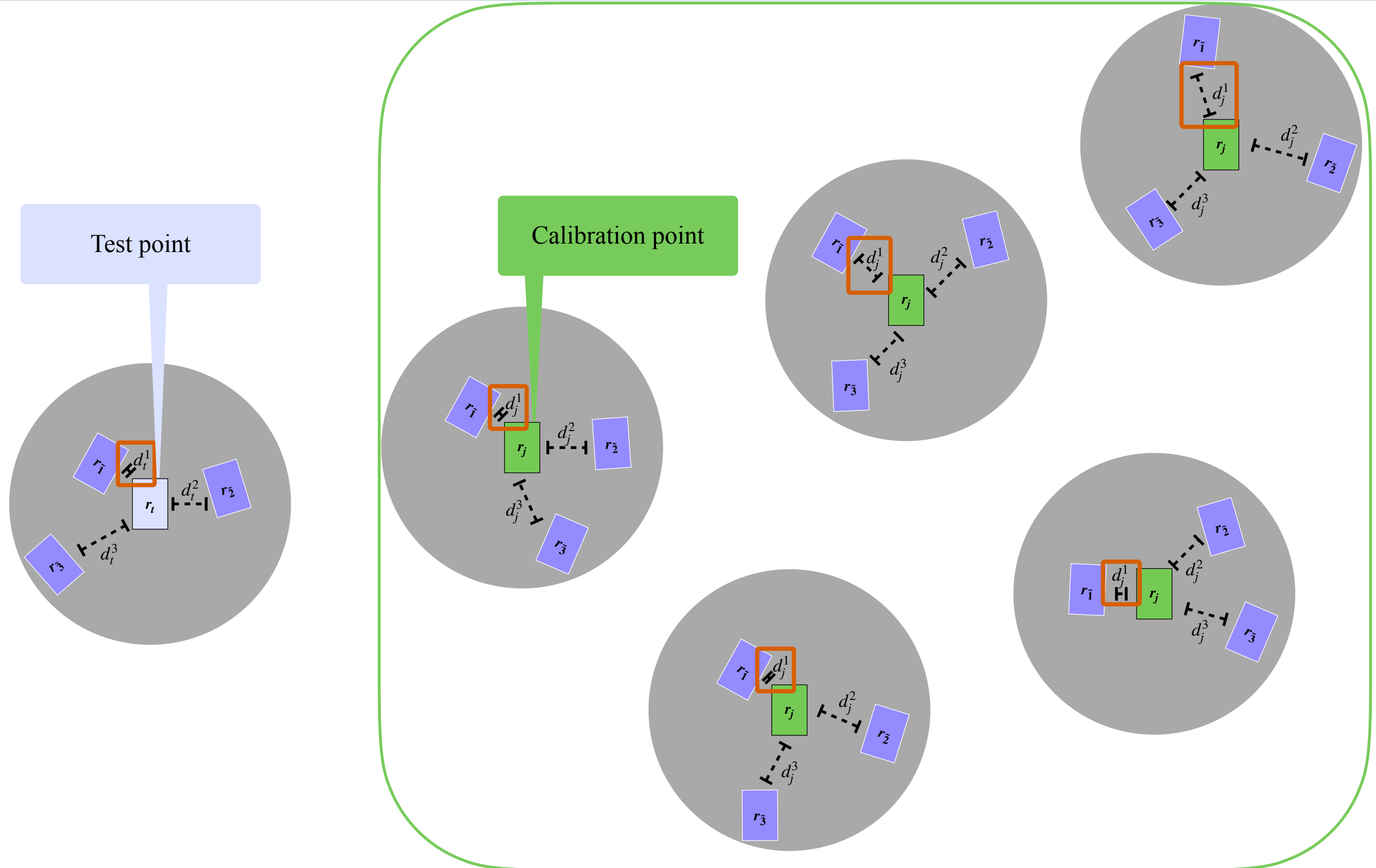


Magnitude of the K-NN Output

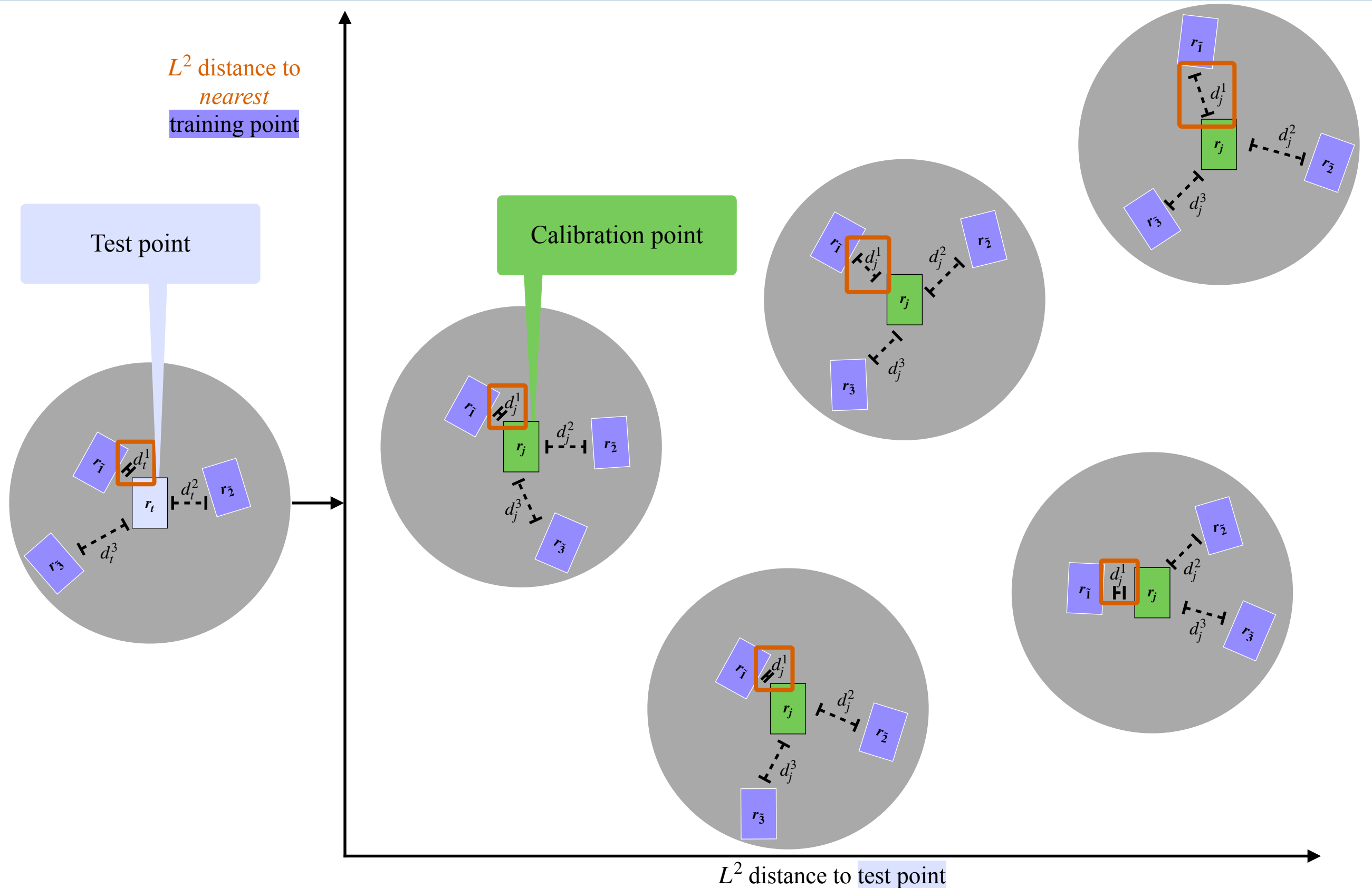
Construct approximation for test point



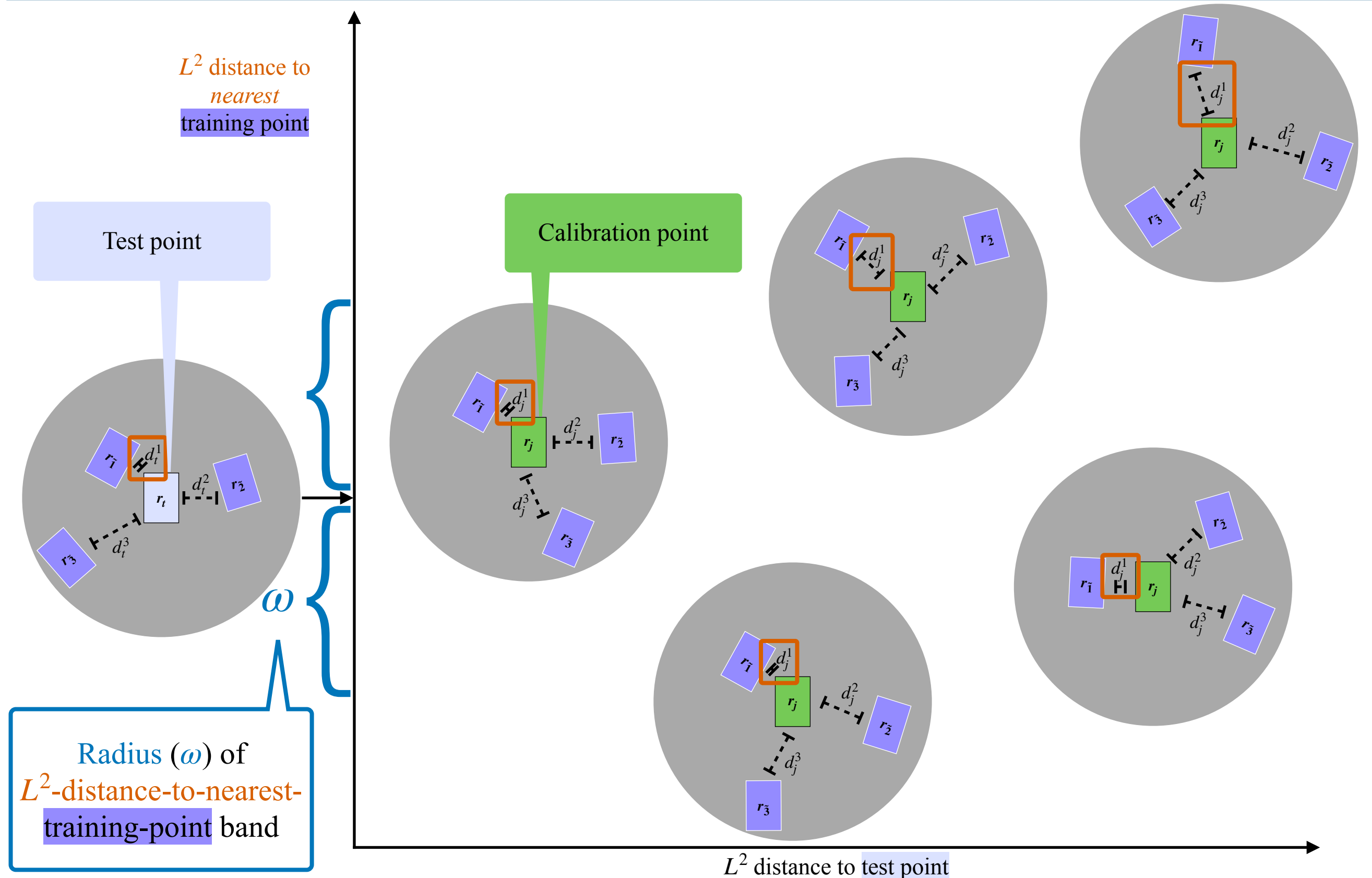
Construct approximation for test point and all calibration points



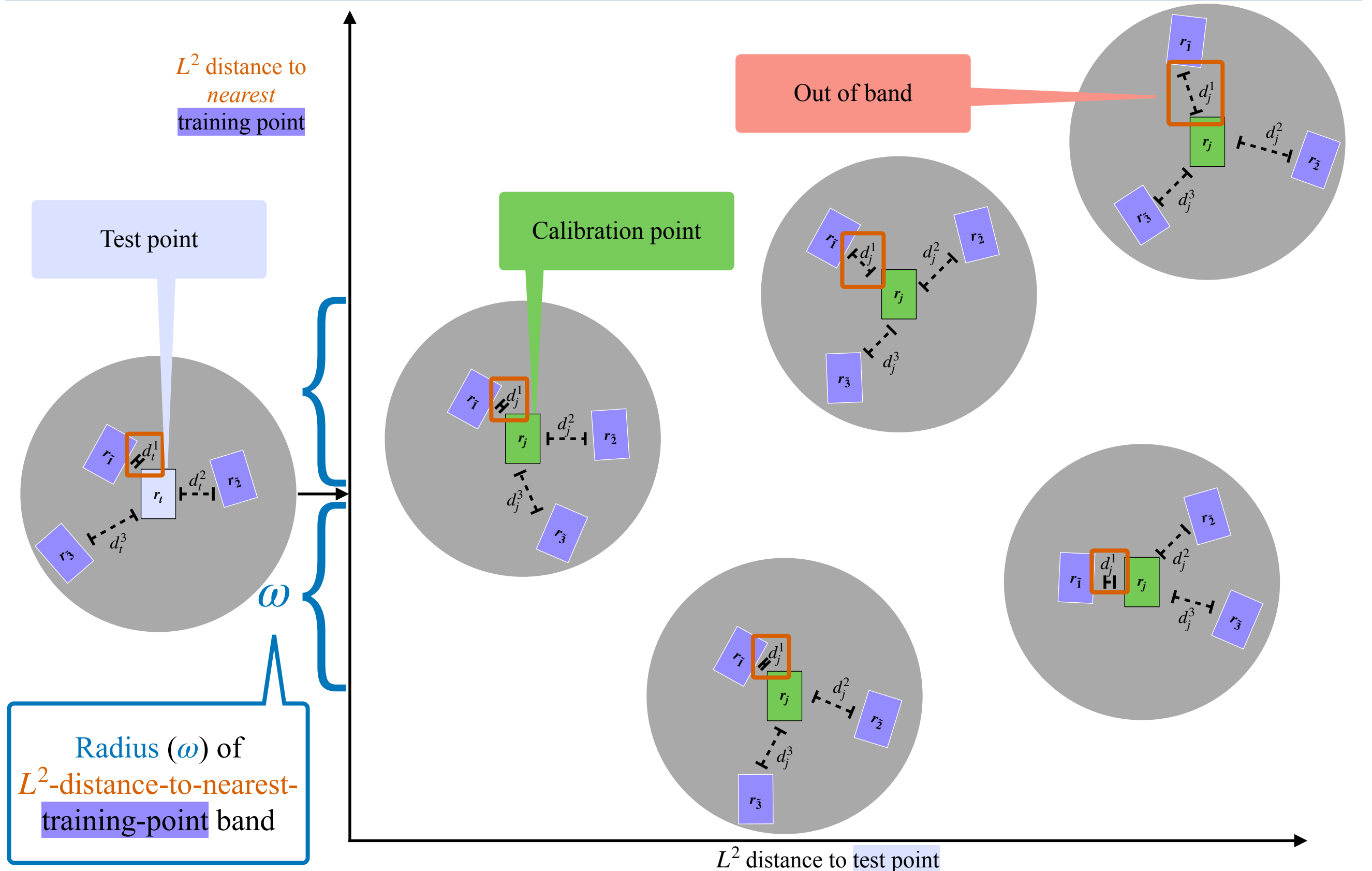
Relate test point to distribution of calibration points



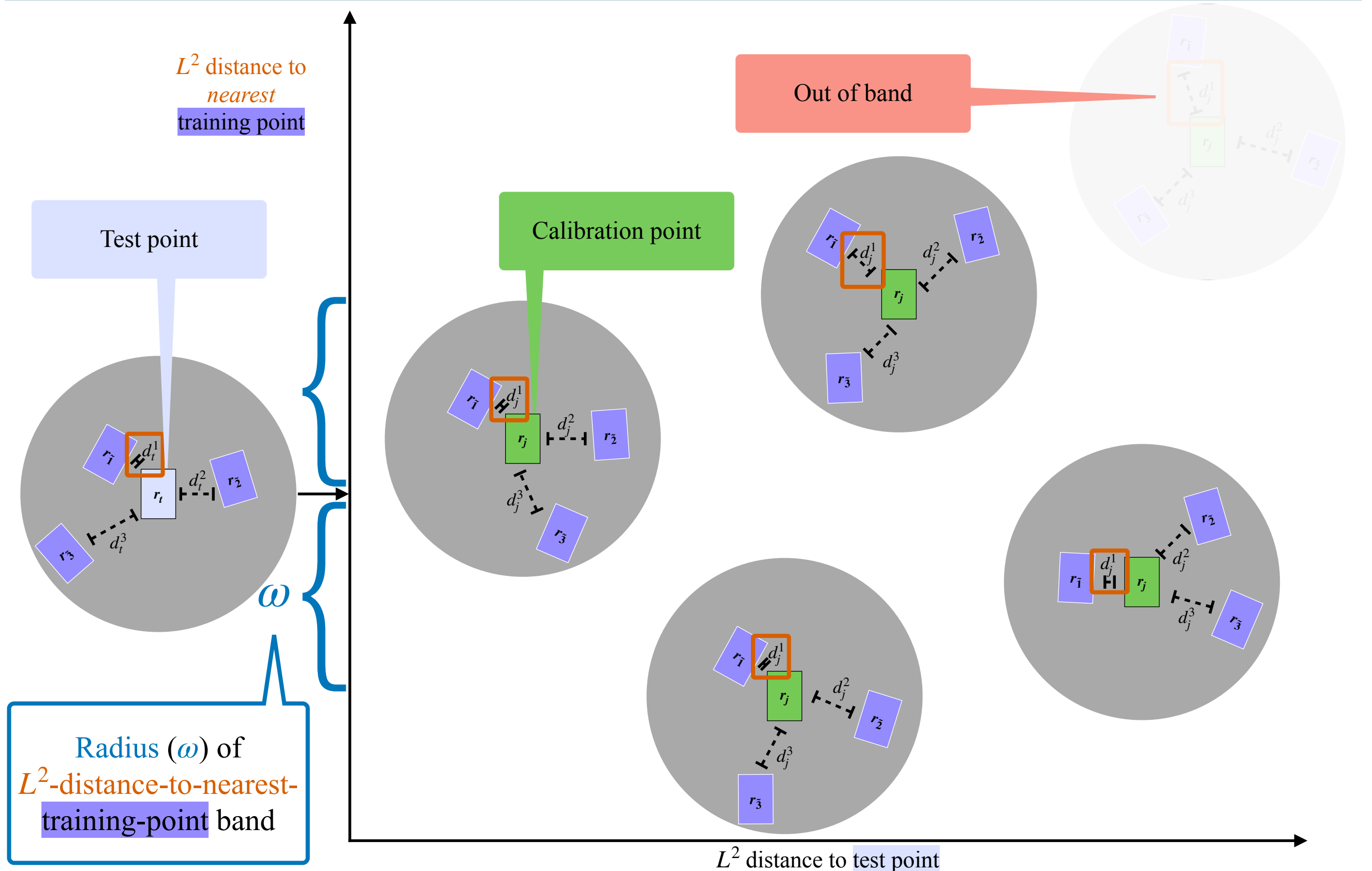
Construct a distance band around the test point



Constrain calibration points to distance band



Constrain calibration points to distance band



Constrain calibration points to distance band

L^2 distance to
nearest
training point

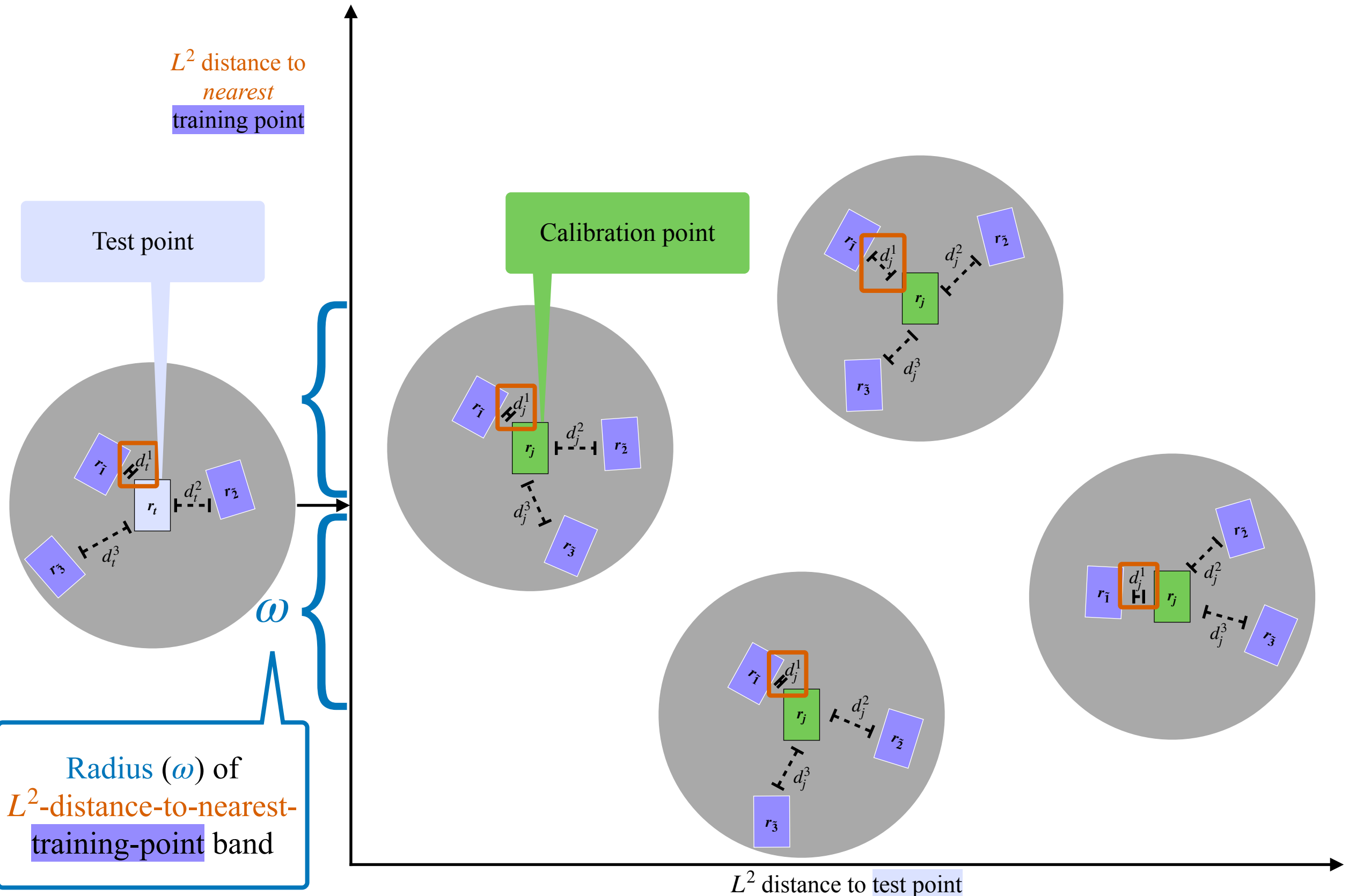
Test point

Calibration point

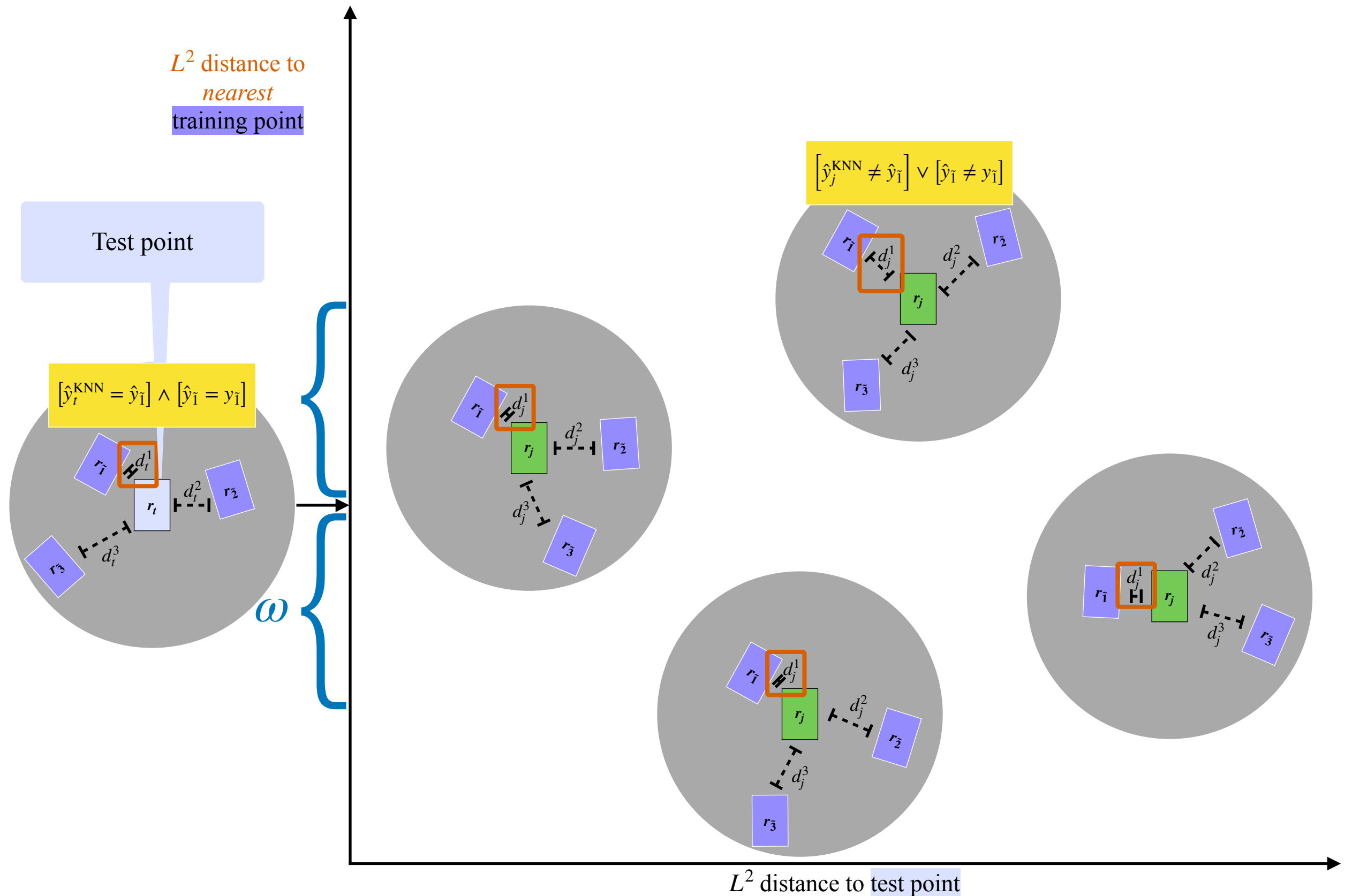
ω

Radius (ω) of
 L^2 -distance-to-nearest-
training-point band

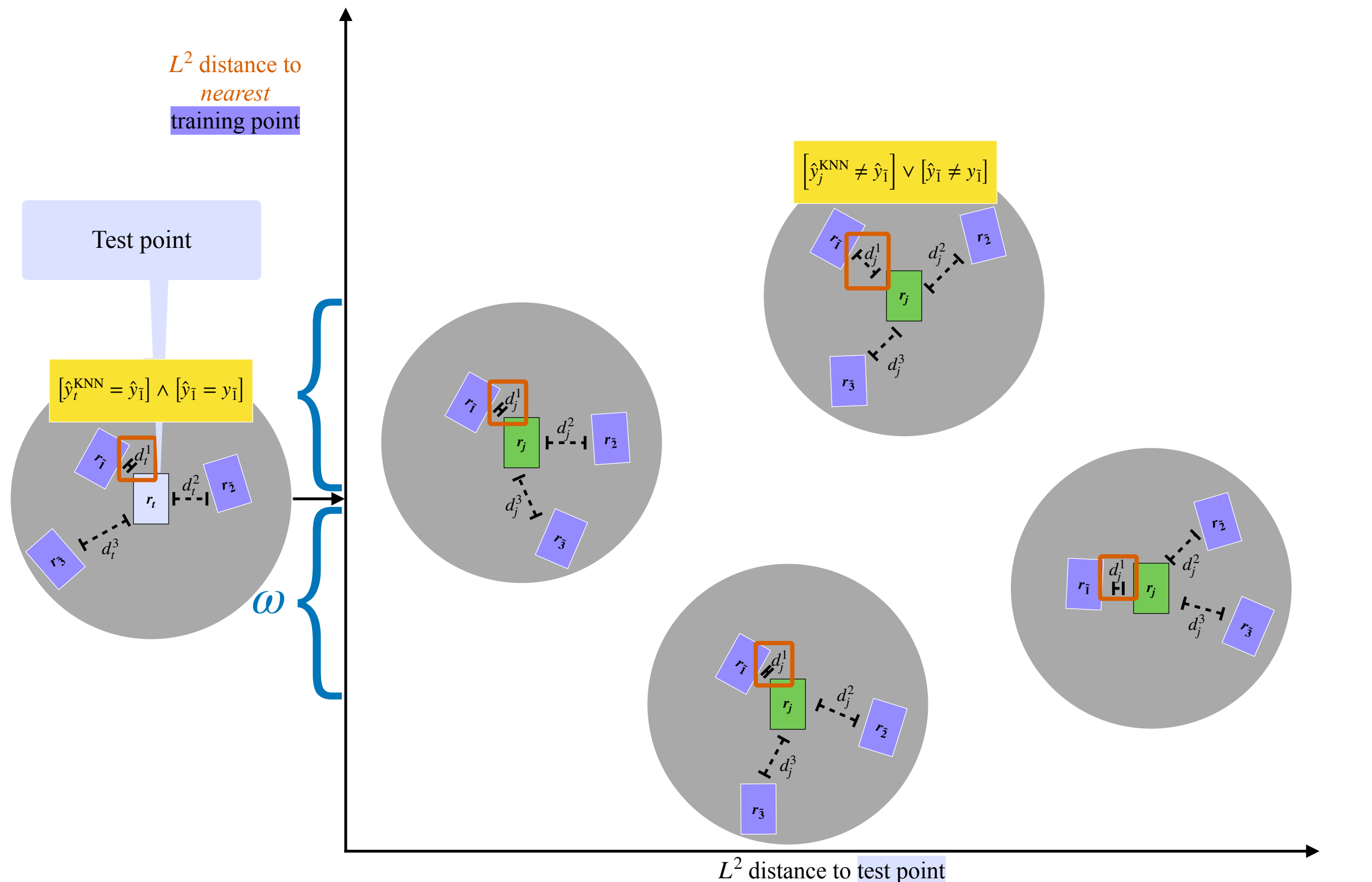
L^2 distance to test point



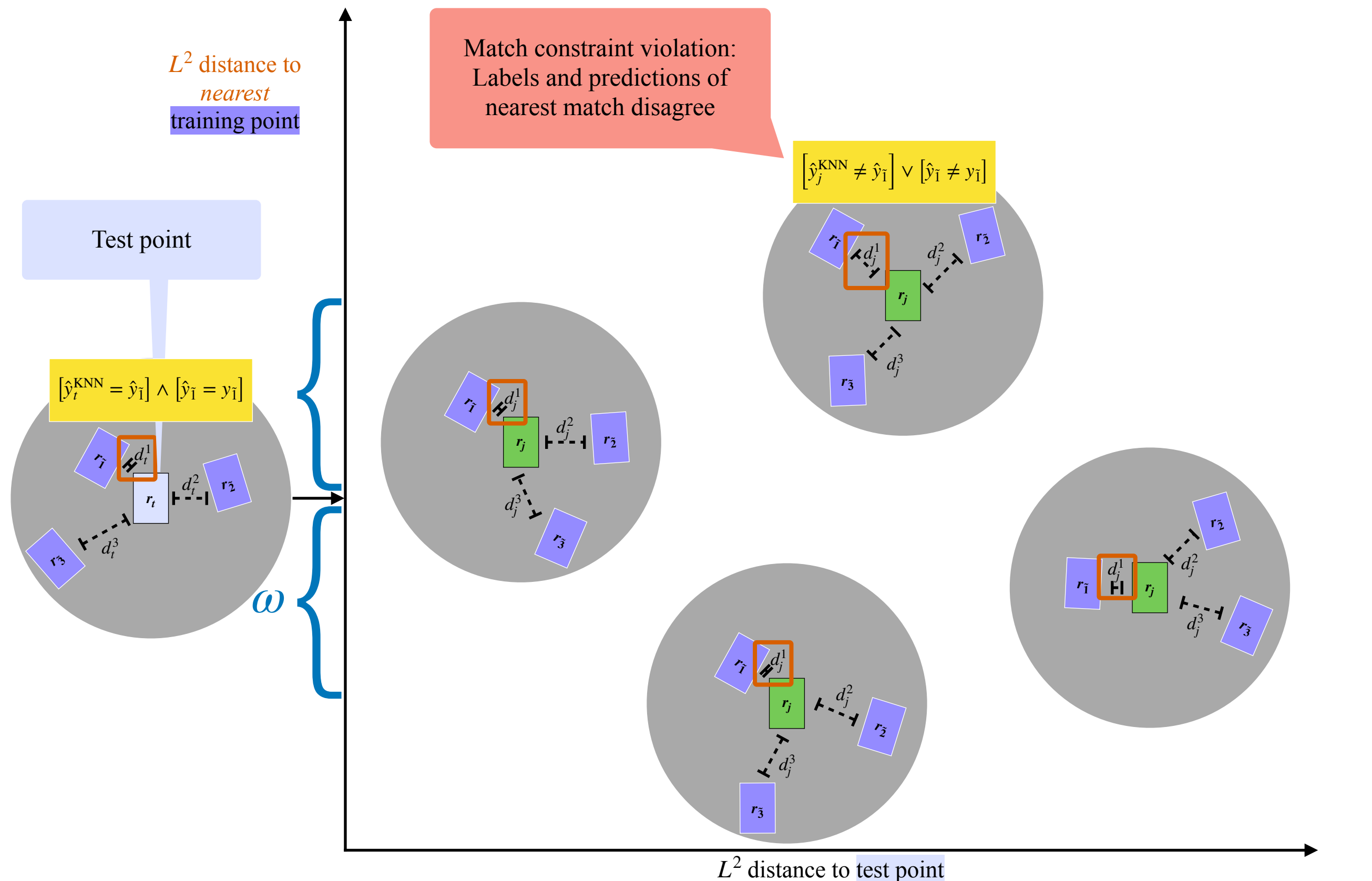
match constraint feature



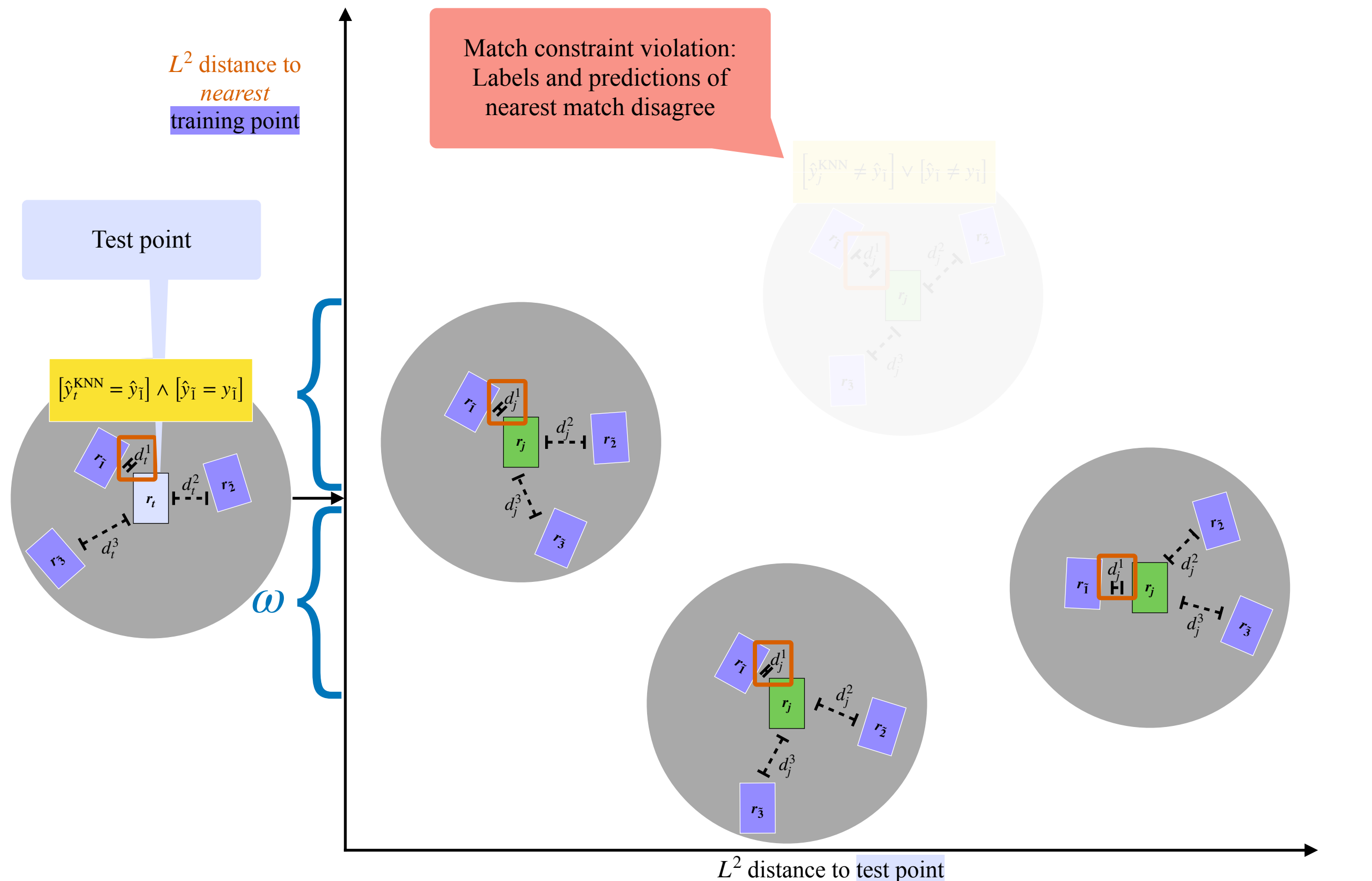
Constrain calibration points to match constraint feature



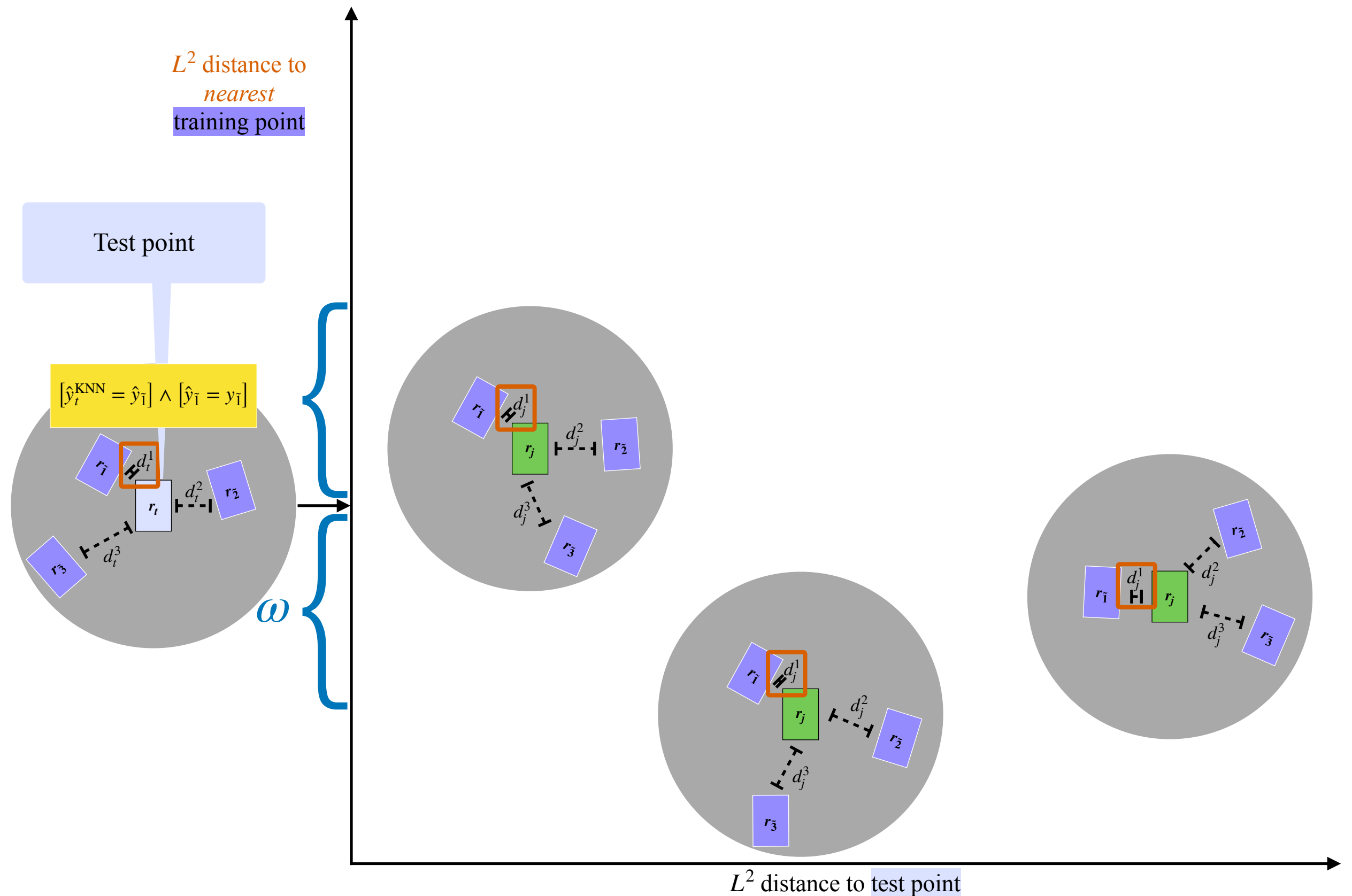
Constrain calibration points to match constraint feature



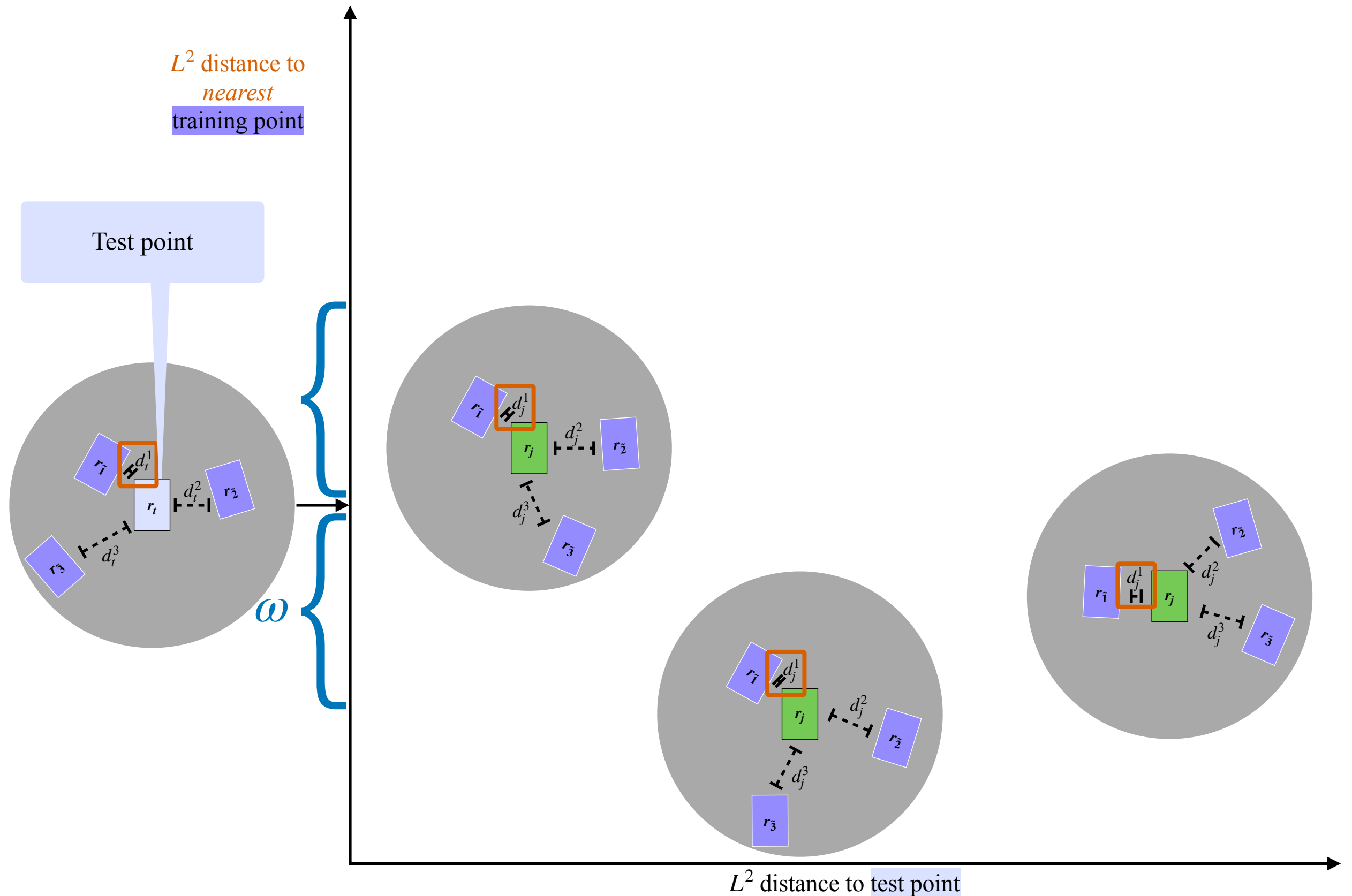
Constrain calibration points to match constraint feature



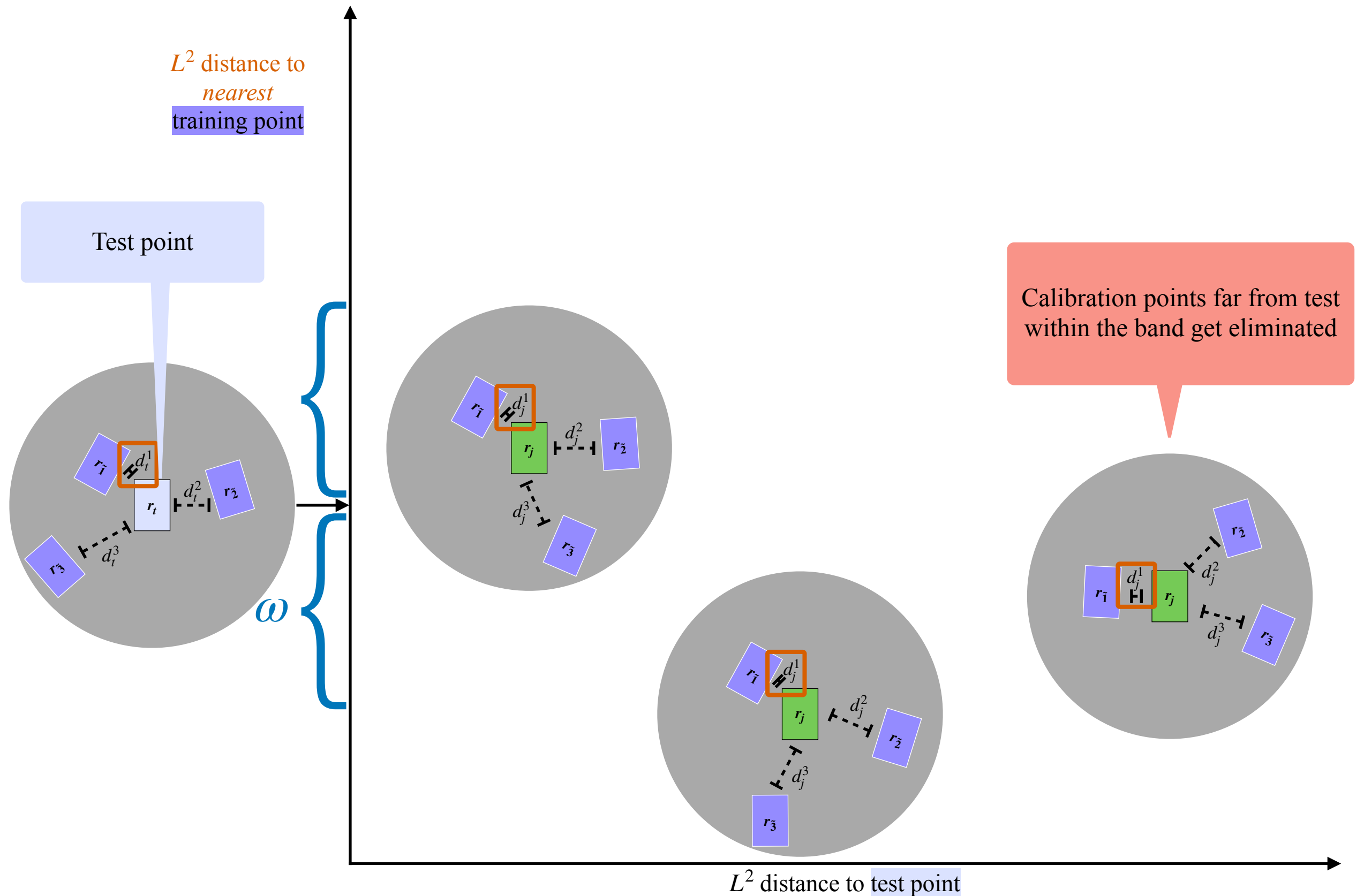
Constrain calibration points to match constraint feature



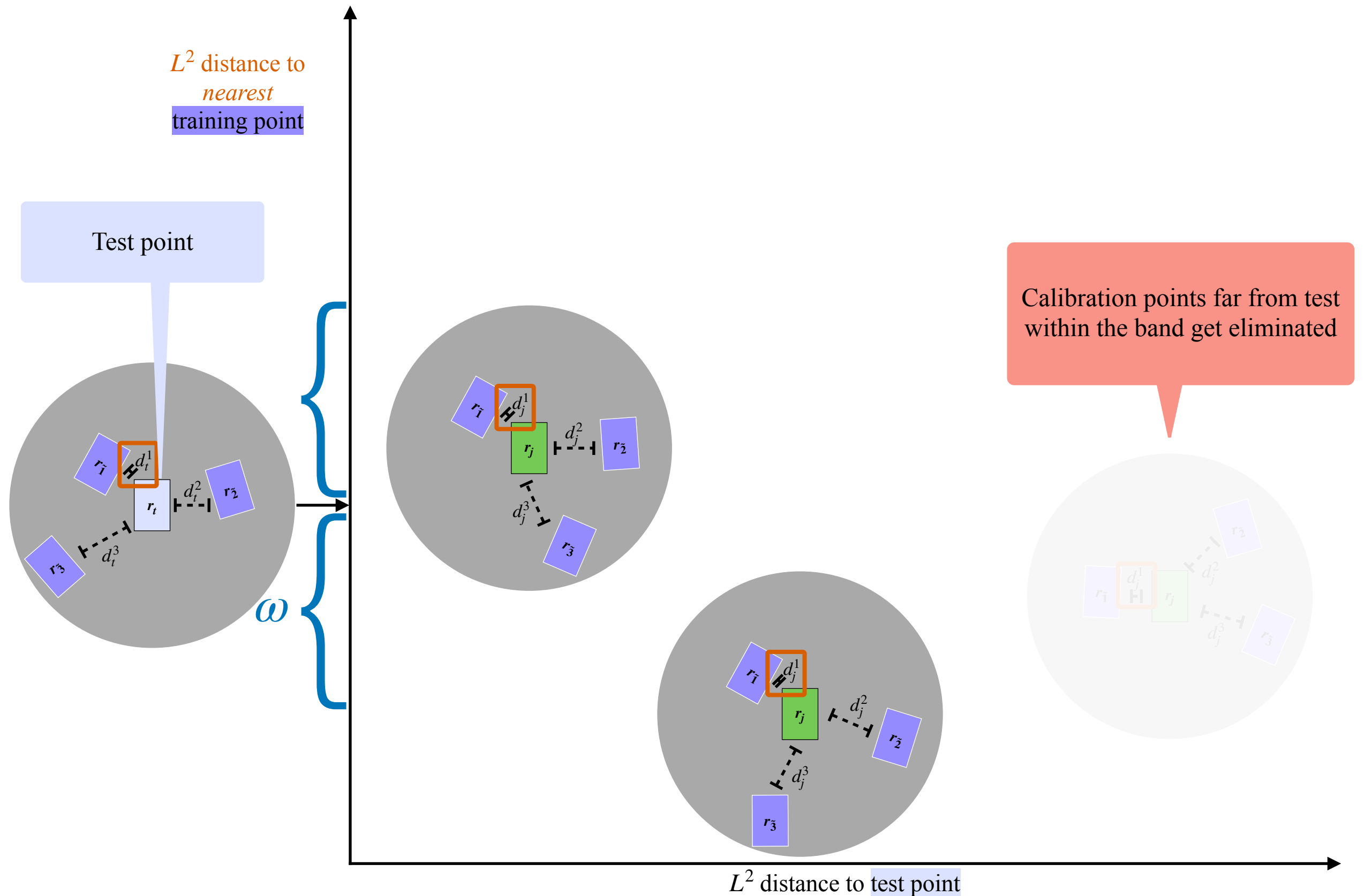
Re-sample \mathcal{D}_{ca} to be more similar to \mathcal{D}_{te}



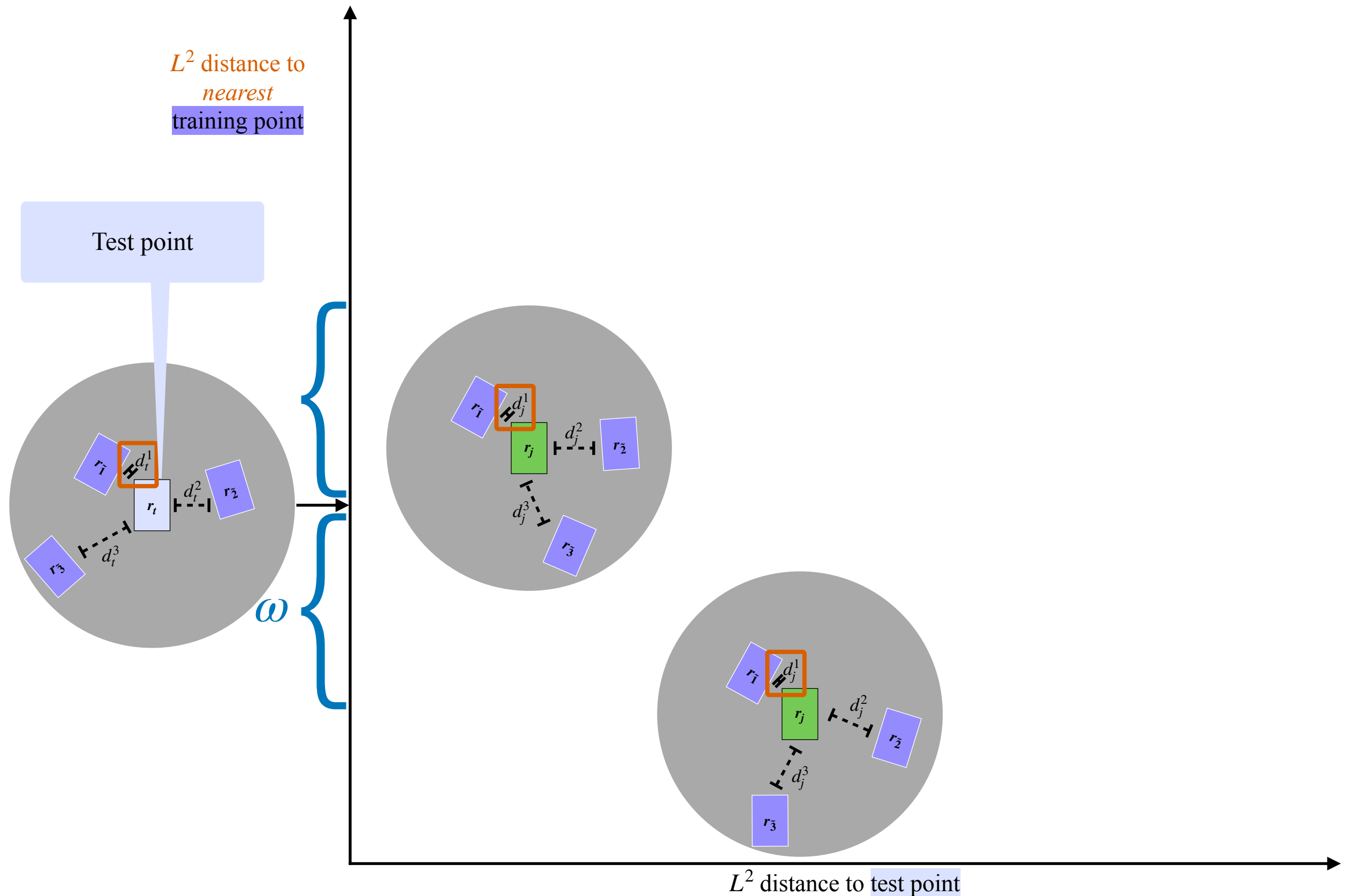
Re-sample \mathcal{D}_{ca} to be more similar to \mathcal{D}_{te}



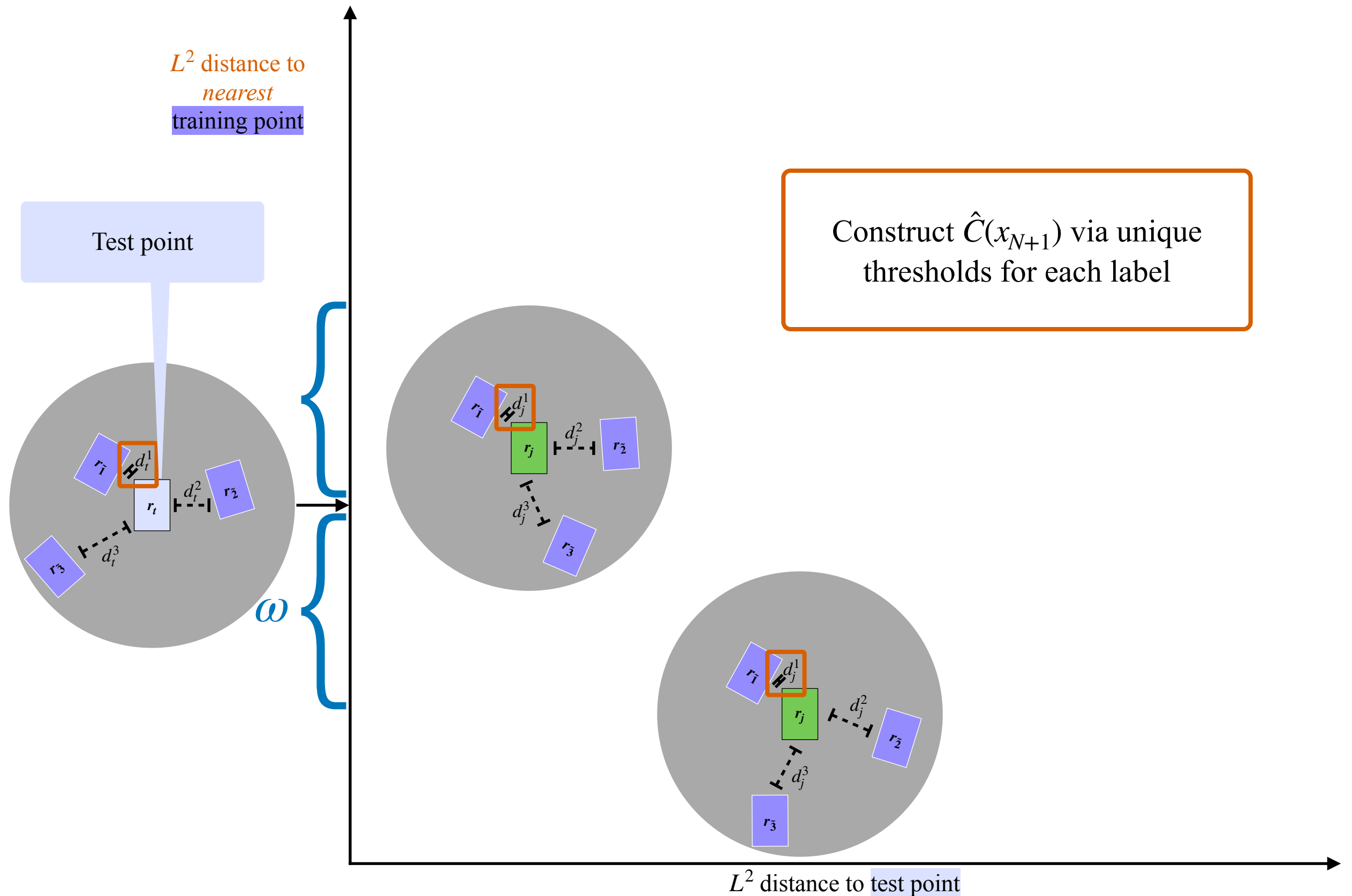
Re-sample \mathcal{D}_{ca} to be more similar to \mathcal{D}_{te}



Re-sample \mathcal{D}_{ca} to be more similar to \mathcal{D}_{te}

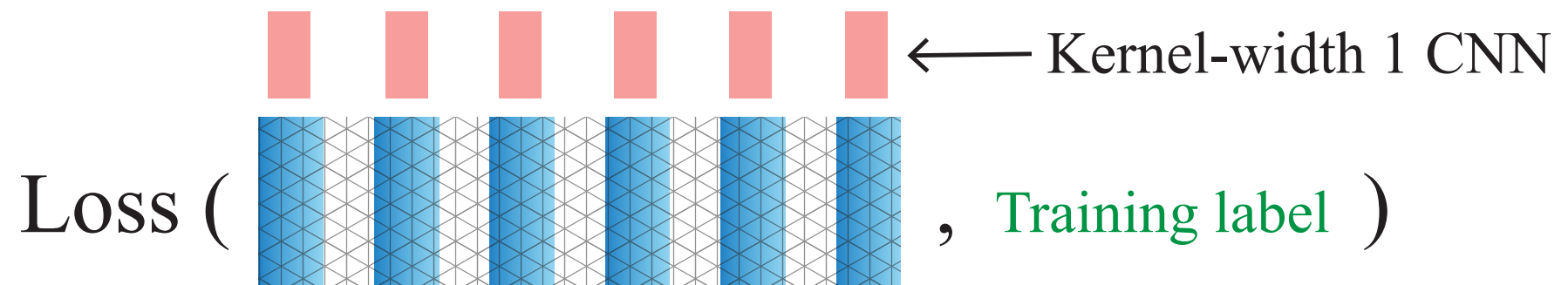


Label-conditional conformal thresholds

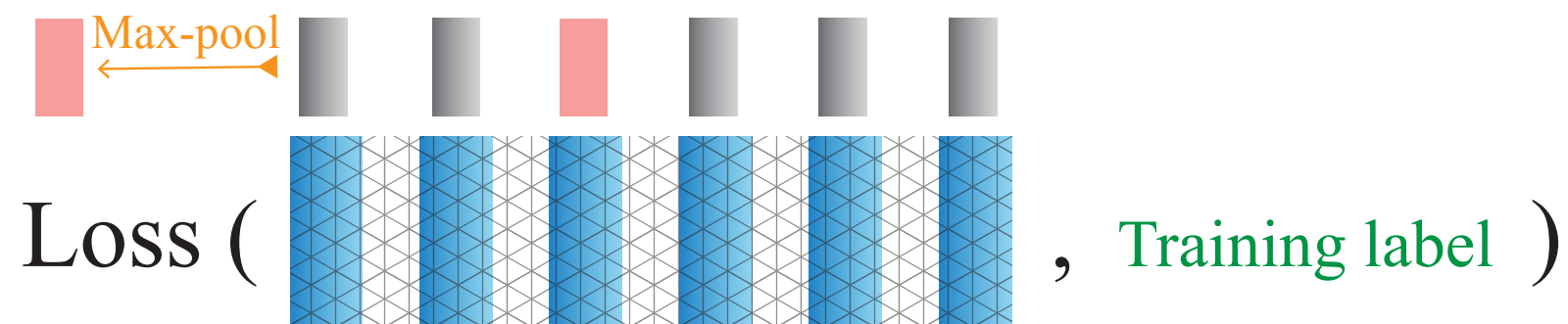


Experiments

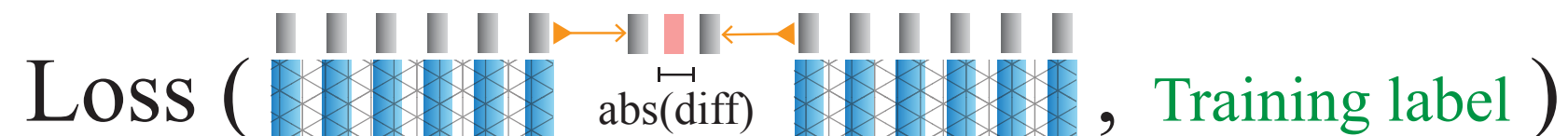
SEQUENCE LABELING:



DOCUMENT CLASSIFICATION (WITH SPARSITY CONSTRAINTS):

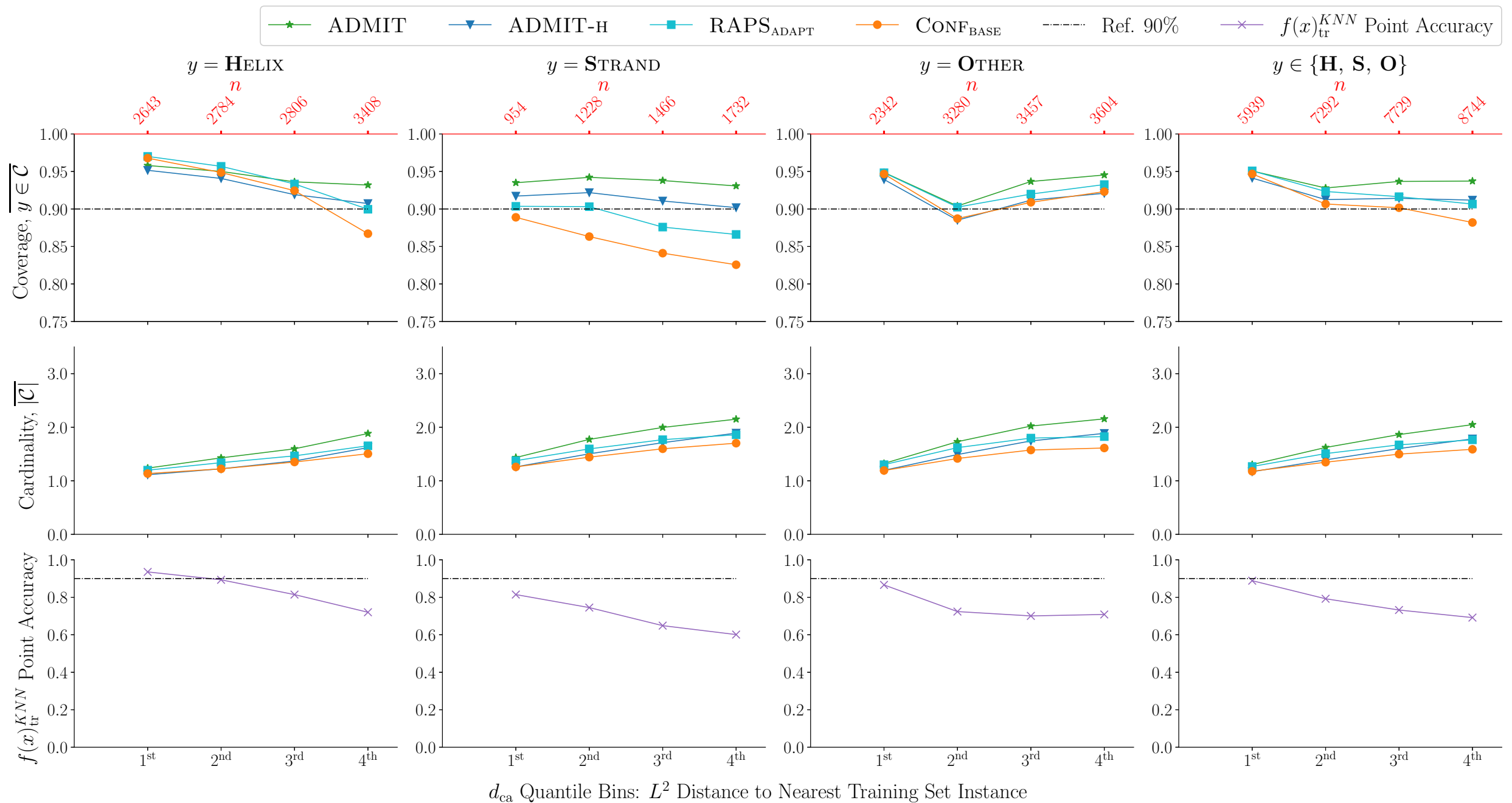


RETRIEVAL-CLASSIFICATION (SEARCH GRAPH):



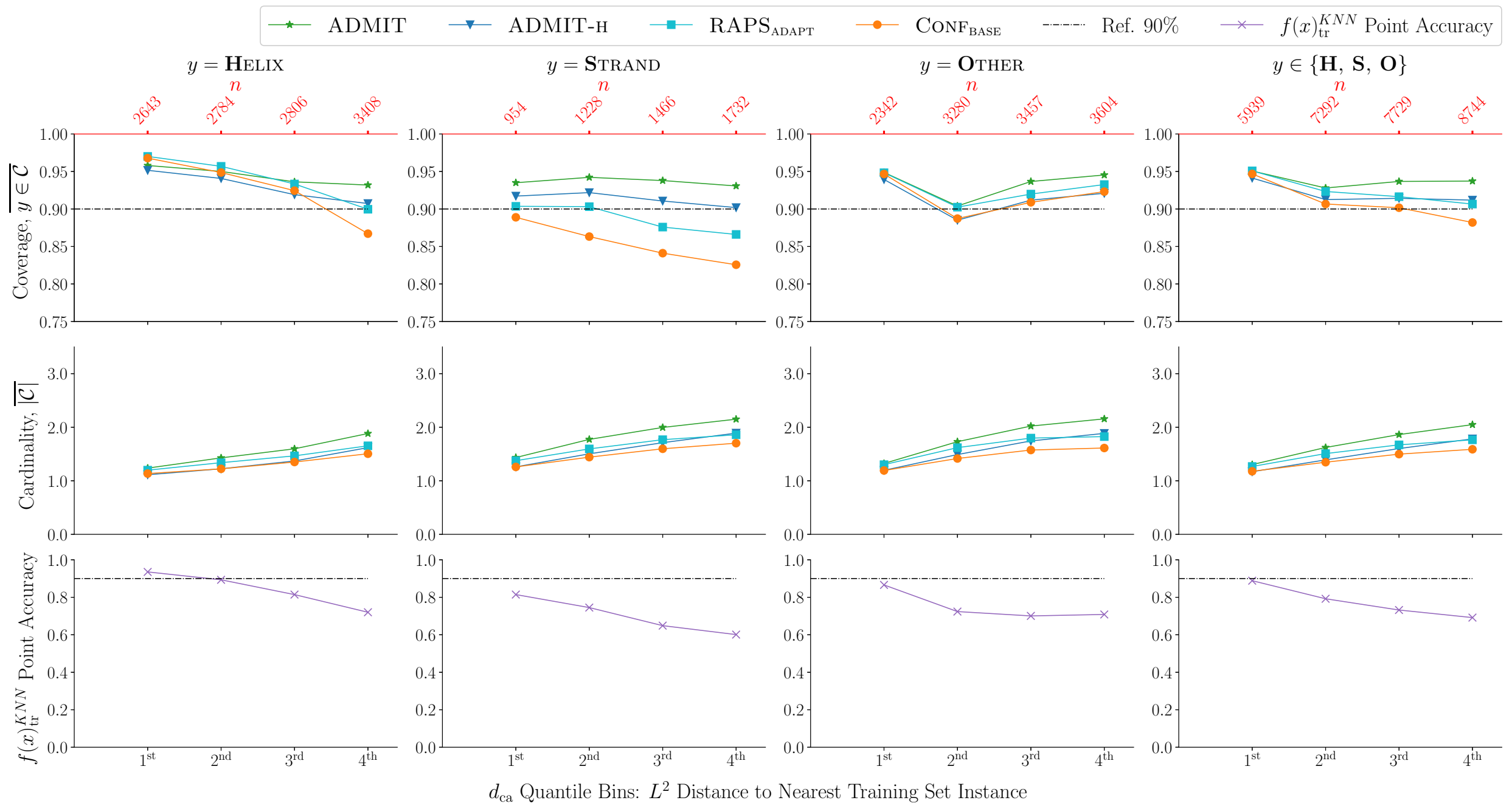
Task diversity: From in-distribution, high-accuracy *to* distribution-shifted, low-accuracy

Empirical behavior on in-domain data



Coverage, cardinality, and point accuracy for the TS115 test set from the PROTEIN task.

Empirical behavior on in-domain data



Coverage, cardinality, and point accuracy for the TS115 test set from the PROTEIN task.

Additional results

- Covariate/label shifts
- Heuristics for coverage conditioned on set composition
- And more ...
- See: <https://arxiv.org/abs/2205.14310>

ADMIT: A general framework for constructing, constraining, and analyzing point predictions and distribution-free prediction sets for deep neural networks.