

# Machine Learning Powered Level-Matching Development Plan

---

## Best Model Choice

**LightGBM Ranker (pairwise ranking) + hard physics masks + Hungarian assignment.**

**Why:** Tree boosting handles nonlinear feature interactions, missing values, mixed discrete and continuous features, and trains fast on small labels. Ranking fits the natural task: for each level in one scheme, rank candidate matches in the other scheme, then enforce one-to-one globally.

---

## End-to-End Pipeline

### 1. Per-Scheme Affine Calibration

Fit robust  $\mathbf{E}' = \mathbf{a}\mathbf{E} + \mathbf{b}$  per dataset against a reference using RANSAC on coarse nearest neighbors. Use residuals in features. Keep  $\mathbf{a}$  and  $\mathbf{b}$  as per-scheme metadata.

### 2. Candidate Generation

For each level  $\mathbf{A}_i$ , collect  $\mathbf{B}_j$  with  $|\mathbf{E}'_i - \mathbf{E}_j| \leq w$  and  $z \leq z_{\max}$  using  $\sigma_{ij} = \sqrt{\sigma_A^2 + \sigma_B^2}$ . Typical  $w = 10$  to  $20$  keV,  $z_{\max} = 4$ .

### 3. Physics Masks

Remove candidates that are hard-forbidden by  $\mathbf{L}$  or  $\mathbf{J}\pi$  rules or by reaction selectivity. Keep a soft prior feature for  $\mathbf{L}$  or  $\mathbf{J}\pi$  compatibility even when not forbidden.

### 4. Features per Pair ( $\mathbf{A}_i, \mathbf{B}_j$ )

- **Core:**  $z, |\Delta E|, \text{sign}(\Delta E)$
- **Quantum numbers:** trinary L match, trinary  $\mathbf{J}\pi$  match, parity match,  $\Delta J$
- **Population priors:** experiment channel flags, known selectivity, beam-target metadata
- **Spectroscopy patterns:**  $\gamma$ -out Jaccard on binned energies, top-k line overlap, intensity similarity, branching vector cosine
- **Structure context:** neighbor spacing similarity before and after, local level density around E, difference in cumulative counts
- **Calibration context:** a and b residuals, per-dataset energy scale uncertainty

*Missing values are fine. LightGBM handles them.*

### 5. Learning Objective

Train **LightGBM Ranker** with pairwise objective. Group by query = A index. Positives are true matches, negatives are other candidates for the same A. If labels are few, start by heuristics for pseudo-labels, then iterate.

## 6. Scoring and Assignment

Predict scores  $\mathbf{s}_{ij}$ . Convert to costs  $\mathbf{C}_{ij} = -\log(\mathbf{s}_{ij} + \epsilon)$ . Keep only candidates that pass hard masks and  $\mathbf{z} \leq \mathbf{z}_{\max}$ . Solve one-to-one with **Hungarian**. If you need monotonicity in energy, use a dynamic-programming matcher with  $i$  increasing implies  $j$  increasing.

## 7. Unmatched Handling

Leave  $\mathbf{A}_i$  or  $\mathbf{B}_j$  unmatched if the best score  $<$  threshold or cost  $>$  cutoff. Report top-k alternates per  $\mathbf{A}_i$  for curator review.

## How to Use with Your Datasets

1. Use your first dataset as reference. Fit affine  $\mathbf{a}, \mathbf{b}$  for each of the other 9 datasets.
2. Generate candidates between the reference and each dataset using your  $\mathbf{E}$  and  $\sigma\mathbf{E}$ , and your  $\mathbf{J}\pi$  and  $\mathbf{L}$  annotations.
3. Start training with a few manually confirmed matches across energy regions. Add hard negatives where  $\mathbf{L}$  or  $\mathbf{J}\pi$  is disallowed.
4. Run prediction to get a ranked match list with top-1 assignment plus top-3 alternates for review.

## Why Not Deep Nets

Tabular, small-to-medium labels, physics rules, and need for interpretability favor gradient-boosted trees and ranking objectives. They are accurate with minimal tuning and give SHAP-style feature attributions for sanity checks.