## oracle_coords_only: Δ PRAUC median

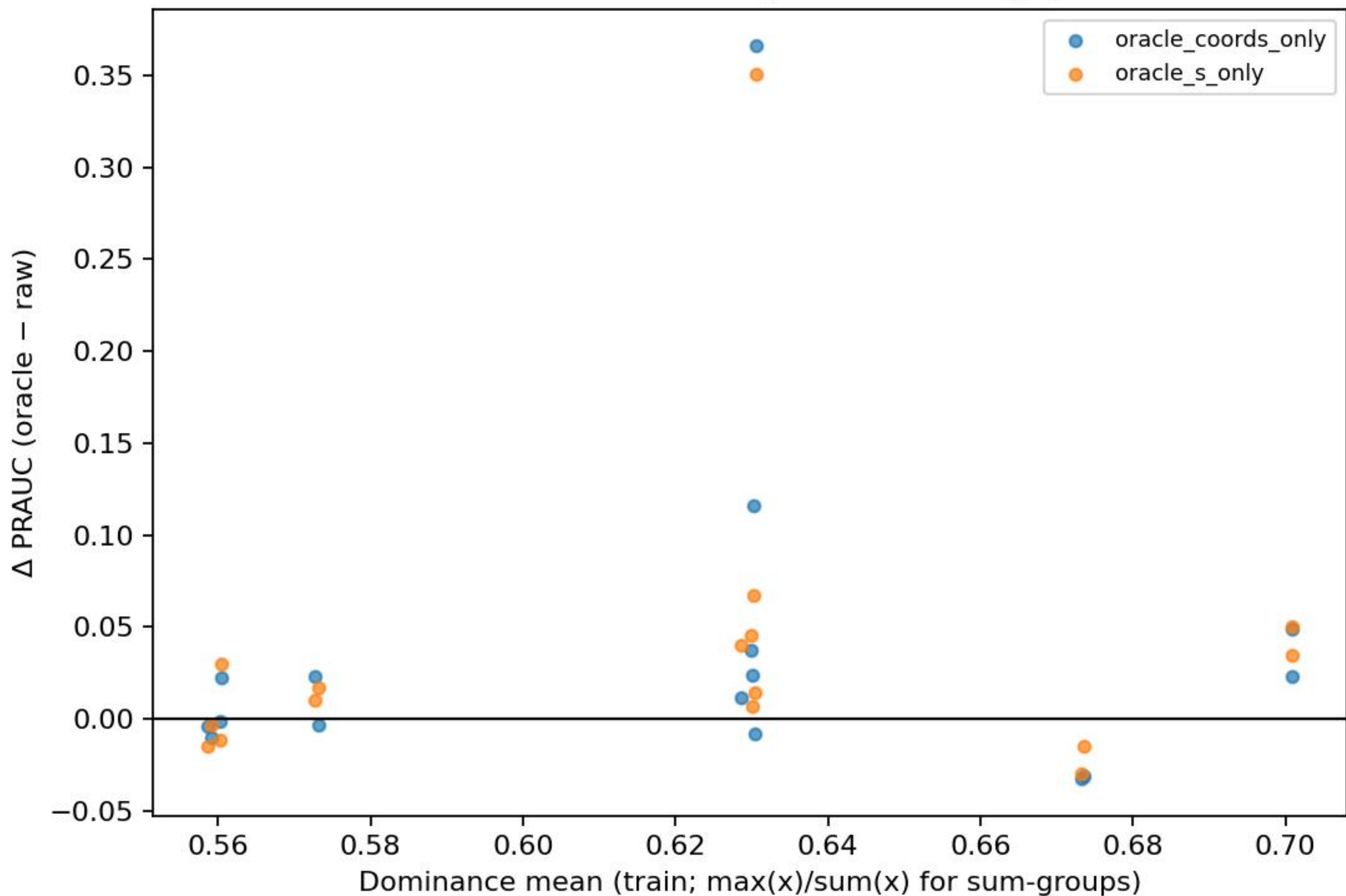|     | corr | mixture_extremes | shift_naive | shift_preserve | tail_heavy |
|-----|------|------------------|-------------|----------------|------------|
| L1  | -0.014 | -0.015 | -0.005 | 0.230 | 0.021 |
| L2  | -0.016 | 0.010 | 0.054 | 0.202 | 0.017 |
| L3  | 0.015 | 0.039 | 0.027 | 0.161 | 0.031 |
| L4  | 0.061 | 0.055 | 0.202 | 0.381 | 0.025 |
| L5  | 0.086 | 0.024 | 0.100 | 0.262 | 0.040 |
| L6  | 0.007 | 0.006 | 0.020 | 0.009 | 0.008 |
| L7  | 0.022 | 0.024 | 0.045 | 0.002 | -0.009 |

## oracle_s_only: Δ PRAUC median

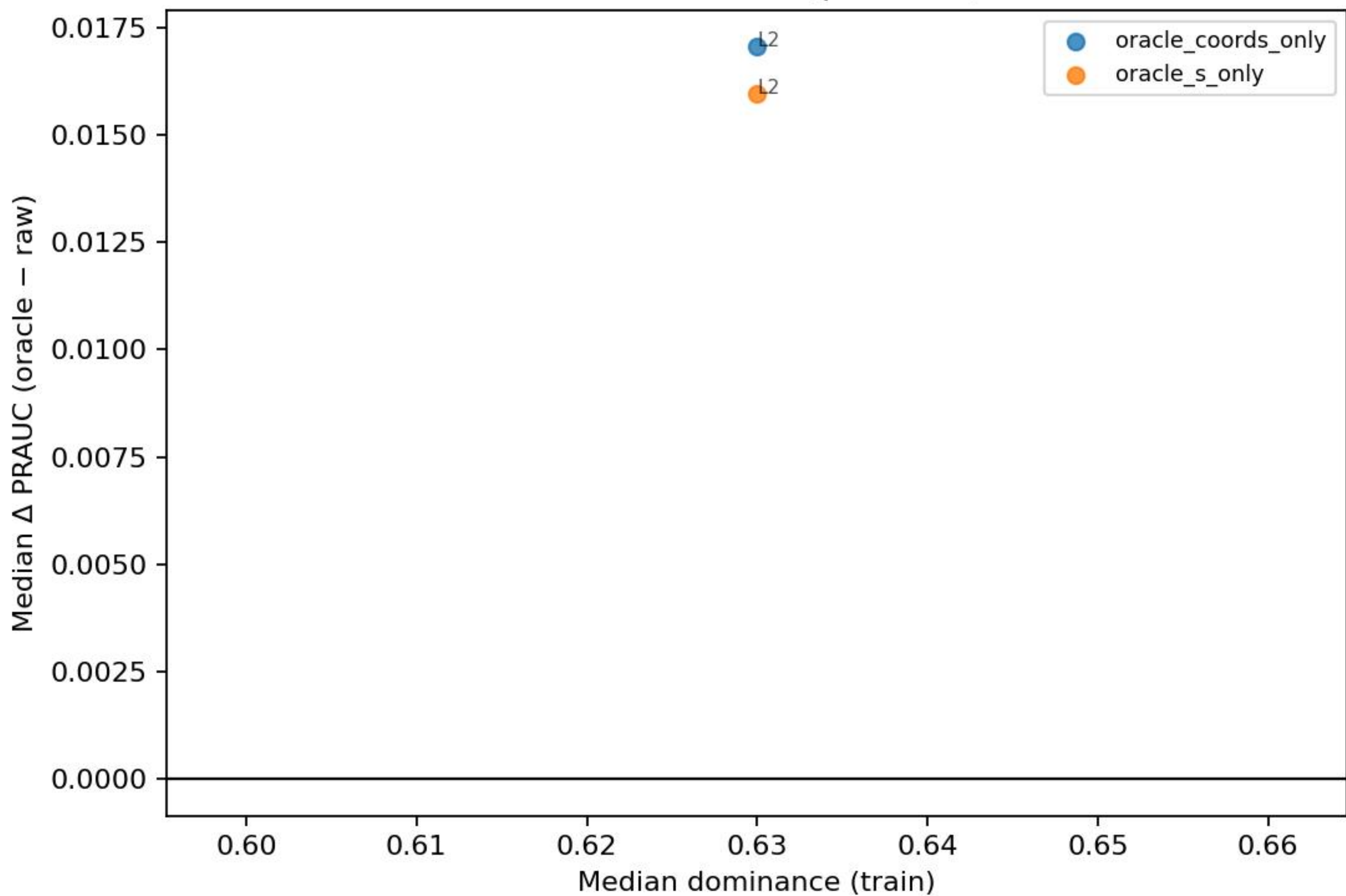|     | corr | mixture_extremes | shift_naive | shift_preserve | tail_heavy |
|-----|------|------------------|-------------|----------------|------------|
| L1  | -0.003 | -0.014 | -0.012 | 0.219 | 0.024 |
| L2  | -0.013 | 0.014 | 0.041 | 0.198 | 0.021 |
| L3  | 0.011 | 0.076 | 0.081 | 0.188 | 0.048 |
| L4  | 0.066 | 0.085 | 0.239 | 0.404 | 0.033 |
| L5  | 0.099 | 0.095 | 0.112 | 0.375 | 0.079 |
| L6  | 0.012 | -0.000 | 0.010 | 0.015 | 0.009 |
| L7  | -0.002 | 0.038 | 0.020 | -0.002 | -0.016 |

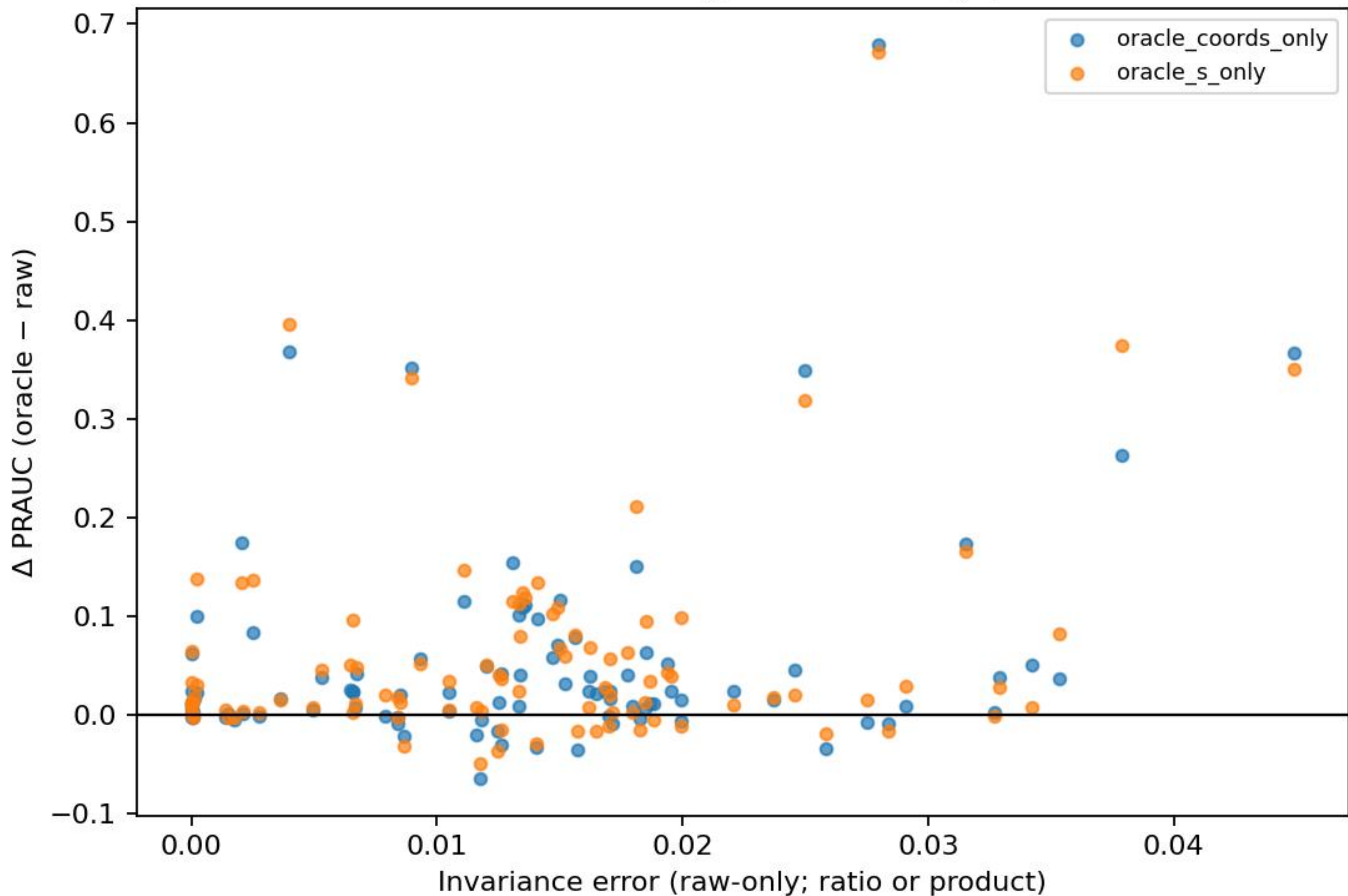Δ PRAUC (oracle − raw) by task and oracle mode
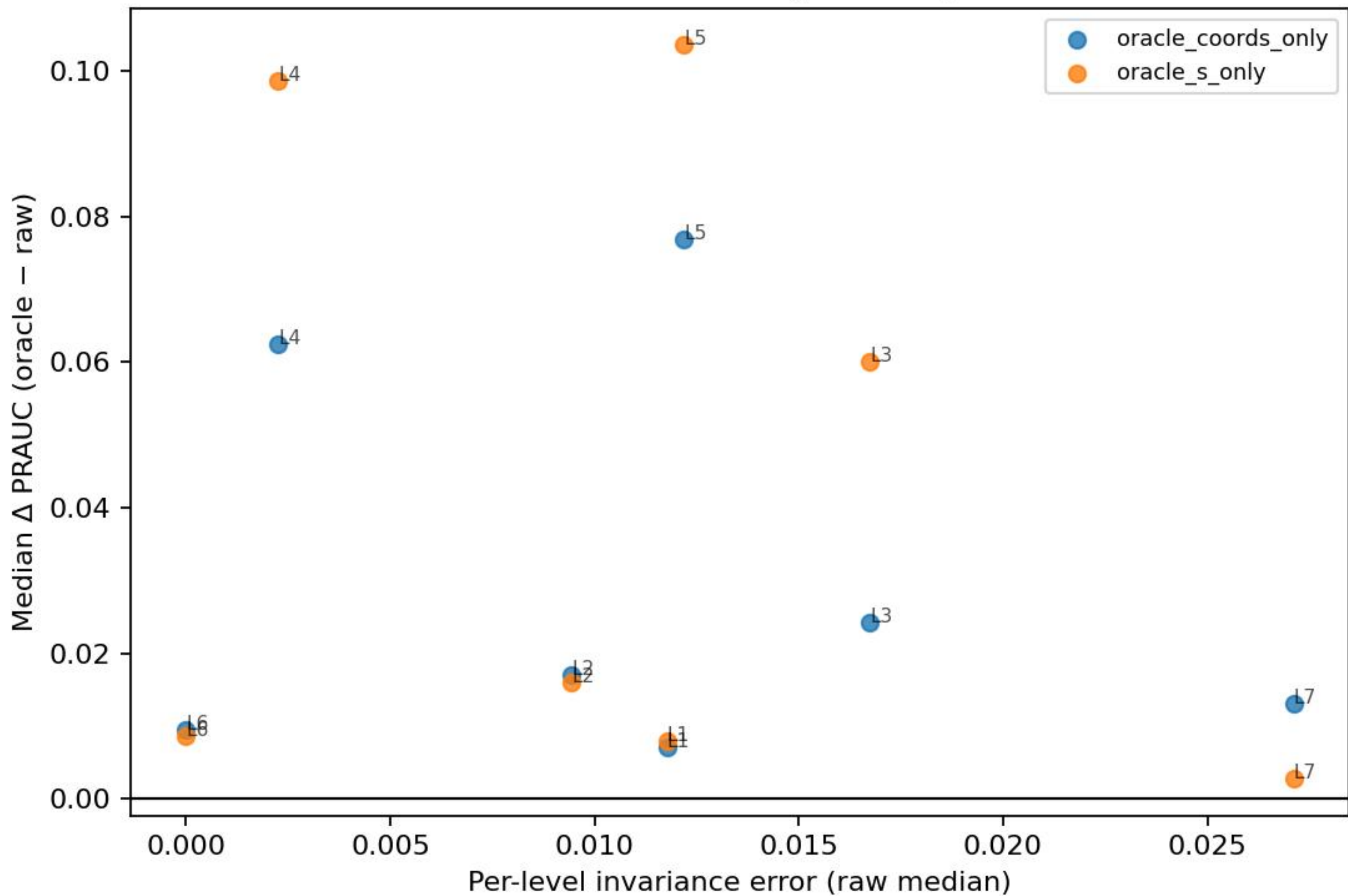
Sum dominance vs generalization gap

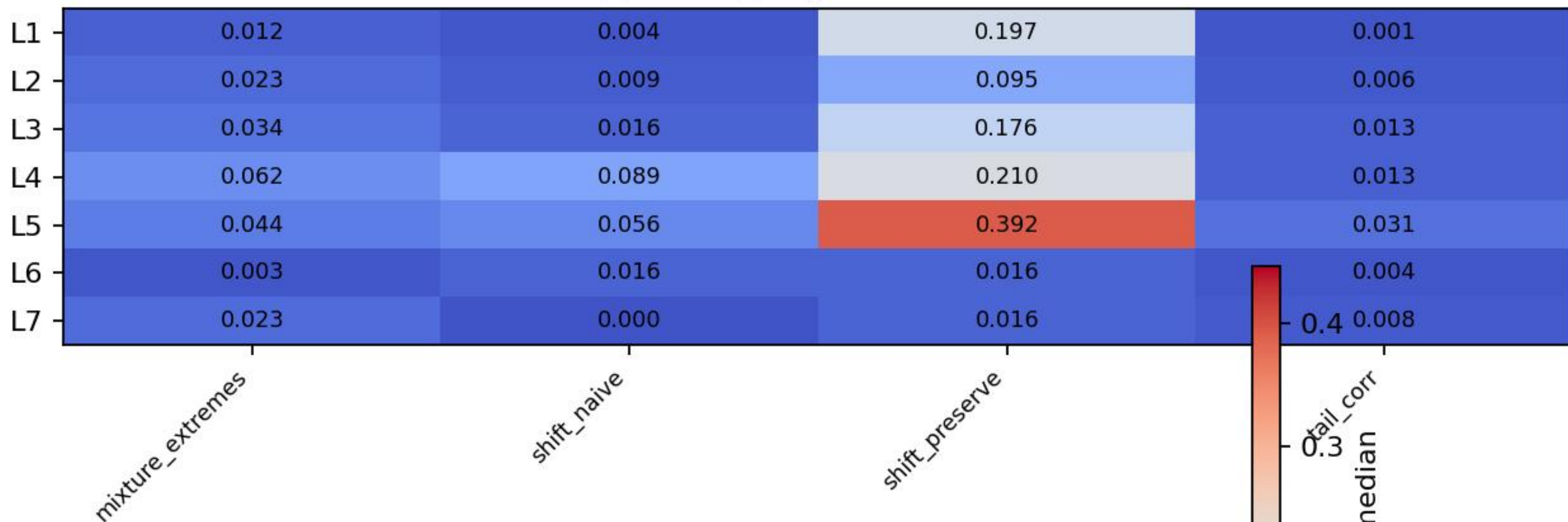Δ vs dominance (per level)

Invariance vs generalization gap

Δ vs invariance (per level)

No n-sweep data

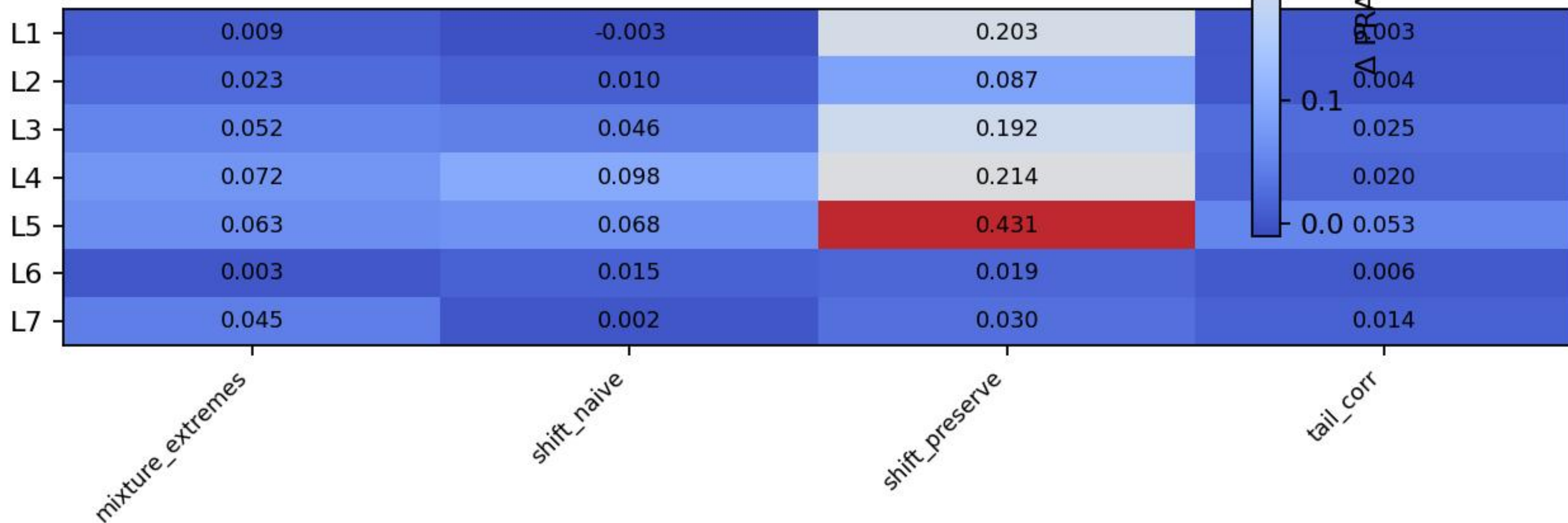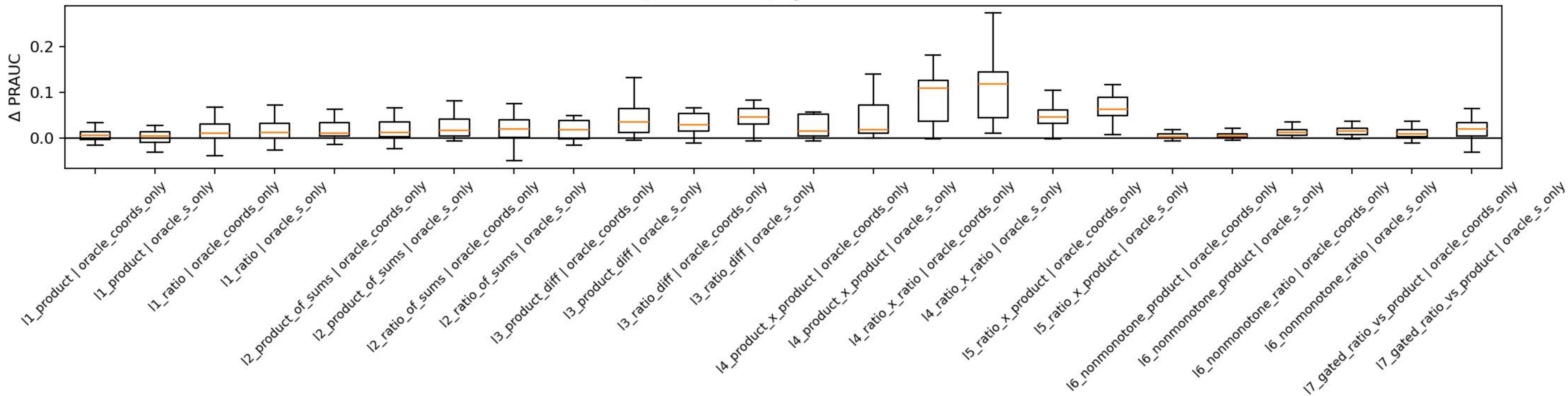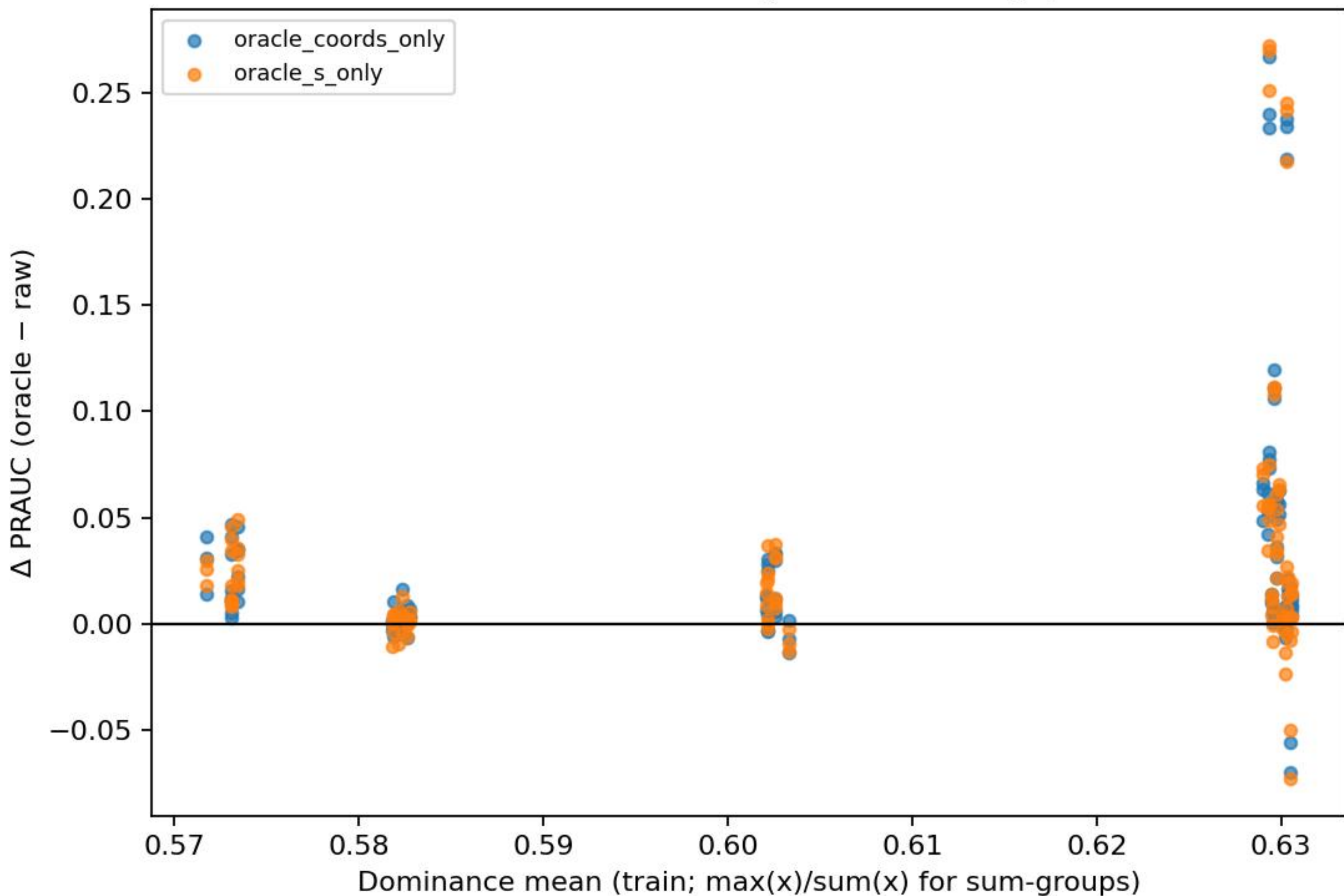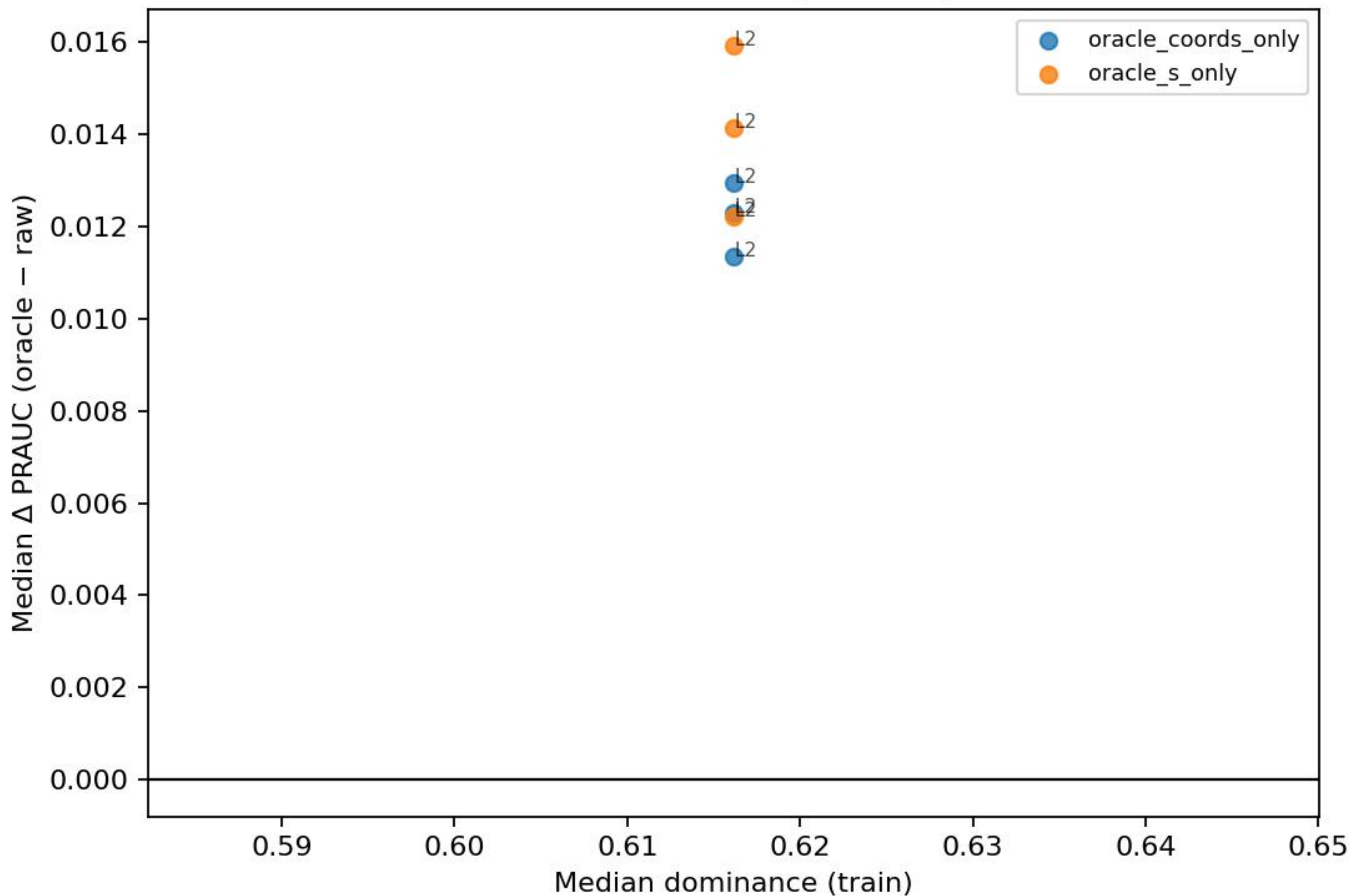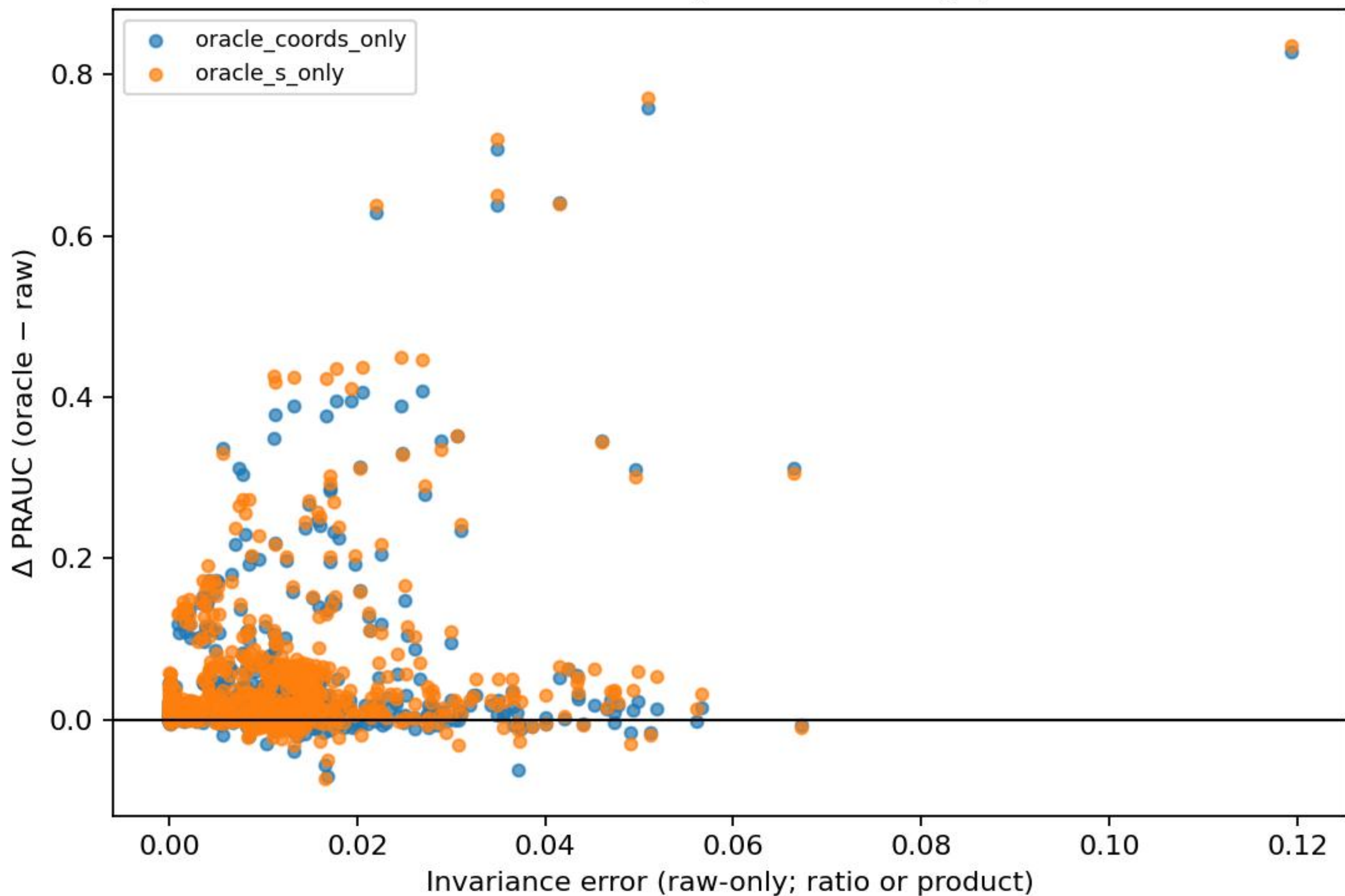**oracle_coords_only: Δ PRAUC median**

| | mixture_extremes | shift_naive | shift_preserve | tail_corr |
|---|---|---|---|---|
| L1 | 0.012 | 0.004 | 0.197 | 0.001 |
| L2 | 0.023 | 0.009 | 0.095 | 0.006 |
| L3 | 0.034 | 0.016 | 0.176 | 0.013 |
| L4 | 0.062 | 0.089 | 0.210 | 0.013 |
| L5 | 0.044 | 0.056 | 0.392 | 0.031 |
| L6 | 0.003 | 0.016 | 0.016 | 0.004 |
| L7 | 0.023 | 0.000 | 0.016 | 0.008 |

**oracle_s_only: Δ PRAUC median**

| | mixture_extremes | shift_naive | shift_preserve | tail_corr |
|---|---|---|---|---|
| L1 | 0.009 | -0.003 | 0.203 | 0.003 |
| L2 | 0.023 | 0.010 | 0.087 | 0.004 |
| L3 | 0.052 | 0.046 | 0.192 | 0.025 |
| L4 | 0.072 | 0.098 | 0.214 | 0.020 |
| L5 | 0.063 | 0.068 | 0.431 | 0.053 |
| L6 | 0.003 | 0.015 | 0.019 | 0.006 |
| L7 | 0.045 | 0.002 | 0.030 | 0.014 |

Δ PRAUC (oracle − raw) by task and oracle mode
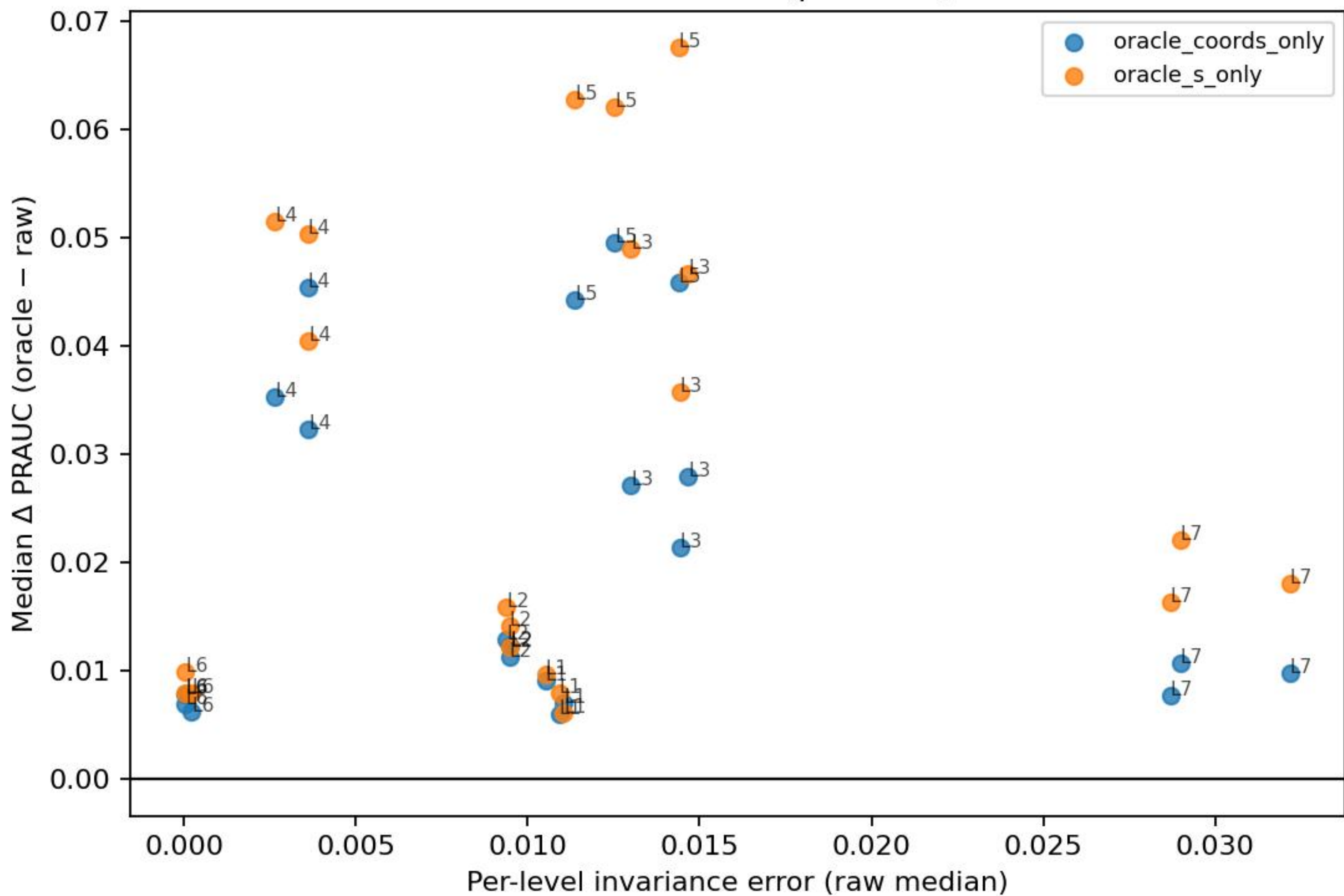
Sum dominance vs generalization gap

Δ vs dominance (per level)

Invariance vs generalization gap

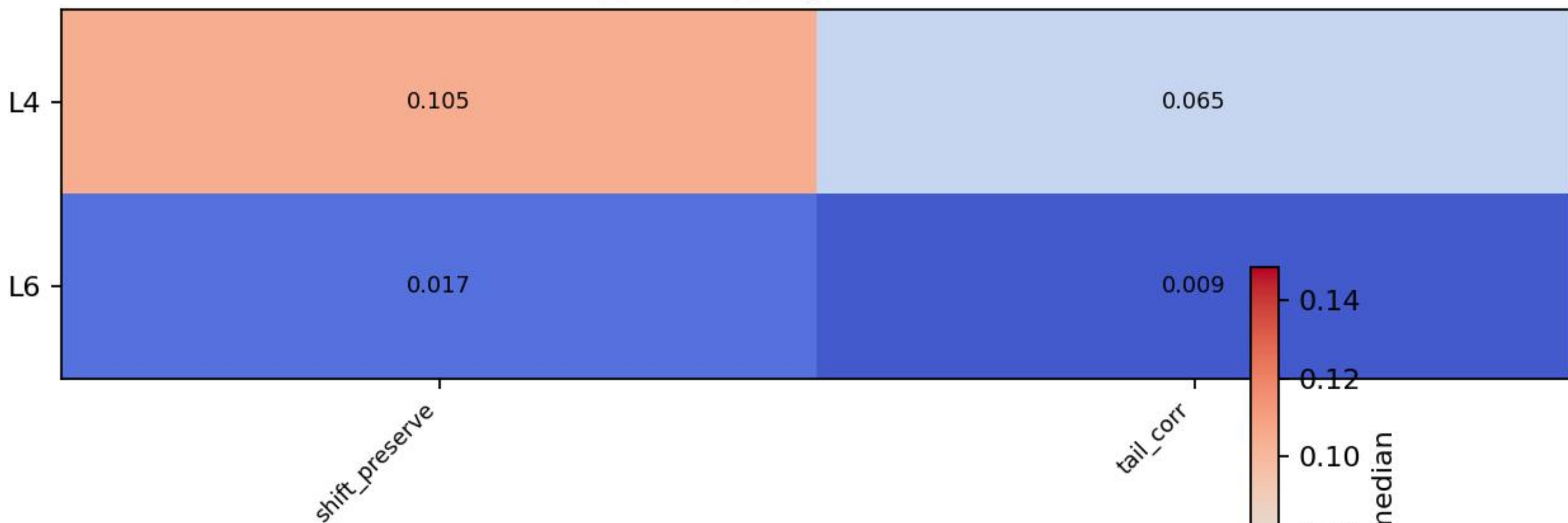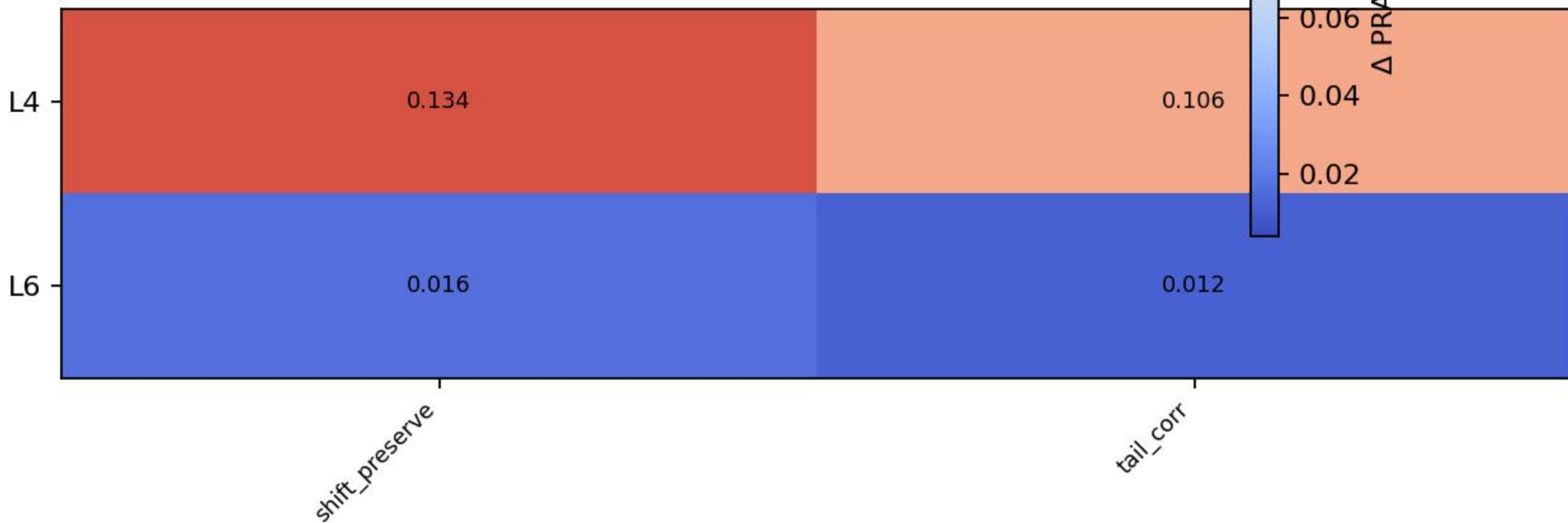Δ vs invariance (per level)

No n-sweep data
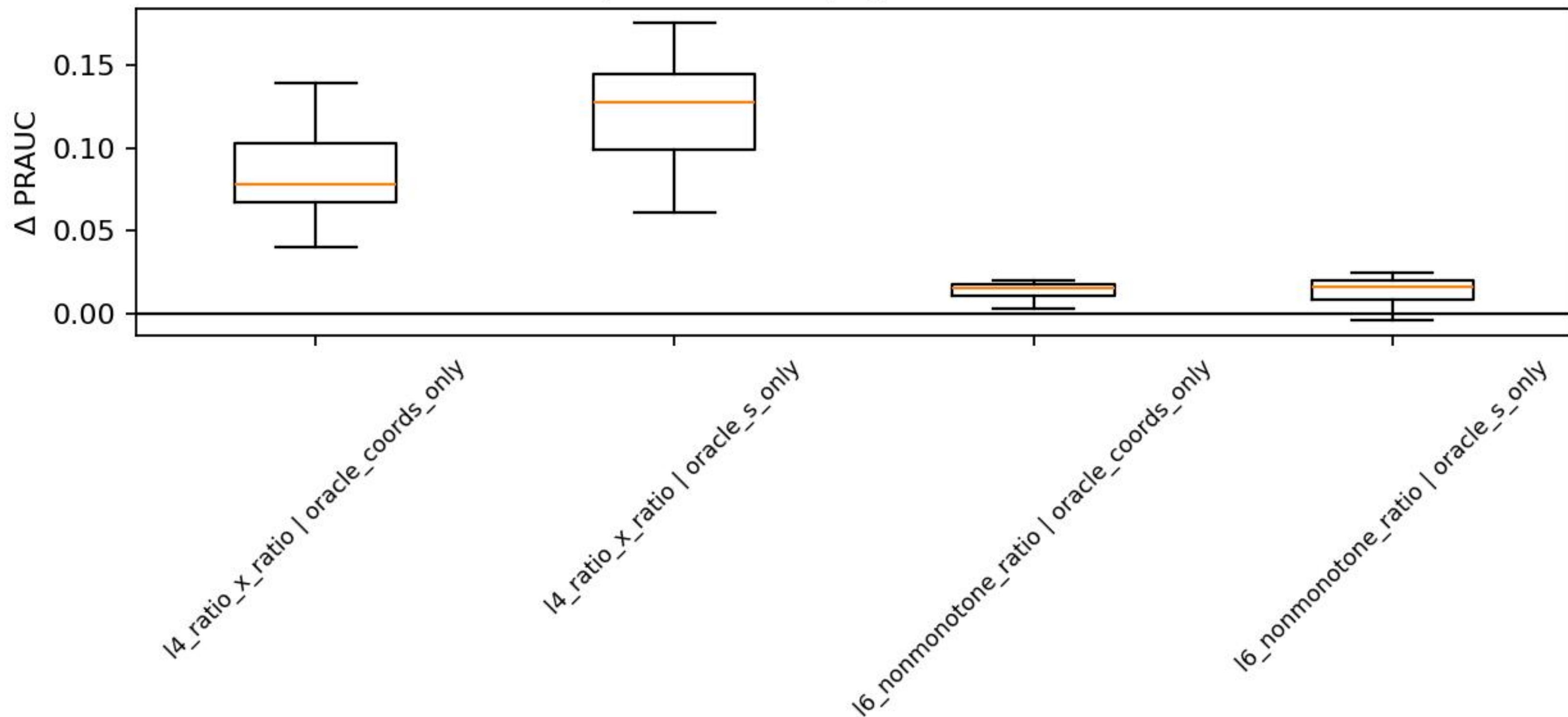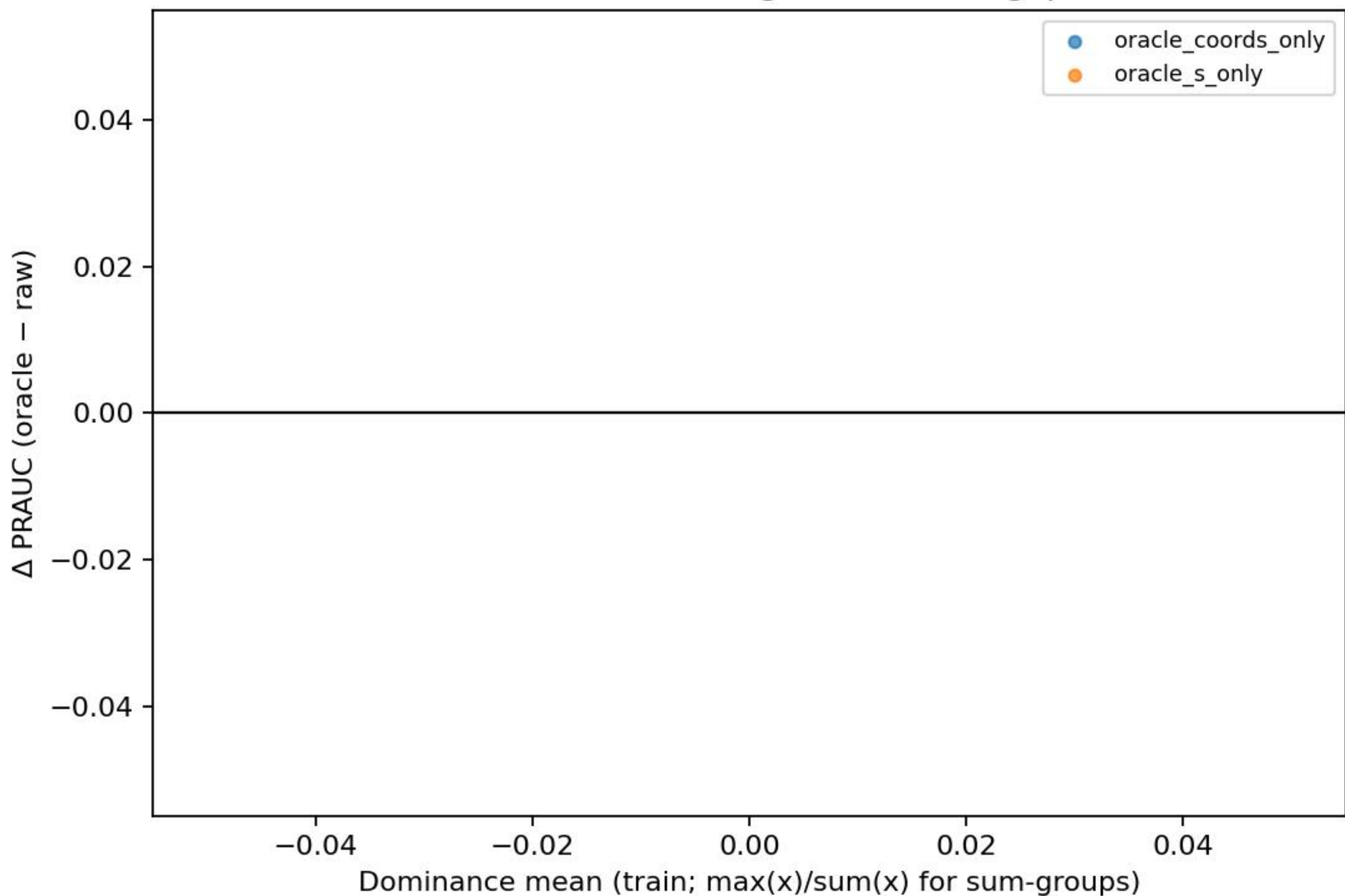
**oracle_coords_only: Δ PRAUC median**

|  | shift_preserve | tail_corr |
|---|---|---|
| L4 | 0.105 | 0.065 |
| L6 | 0.017 | 0.009 |

**oracle_s_only: Δ PRAUC median**

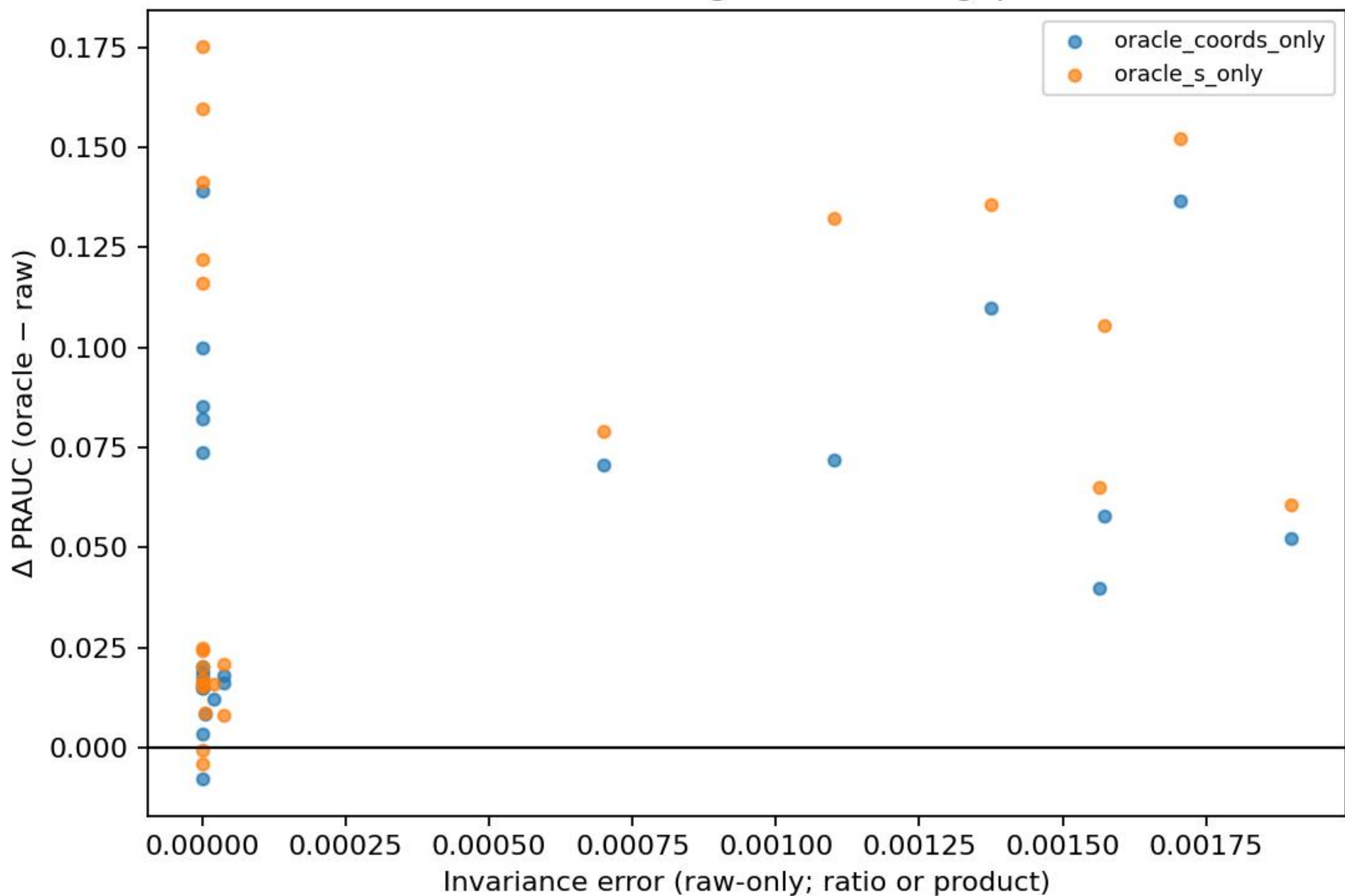|  | shift_preserve | tail_corr |
|---|---|---|
| L4 | 0.134 | 0.106 |
| L6 | 0.016 | 0.012 |

Δ PRAUC (oracle − raw) by task and oracle mode

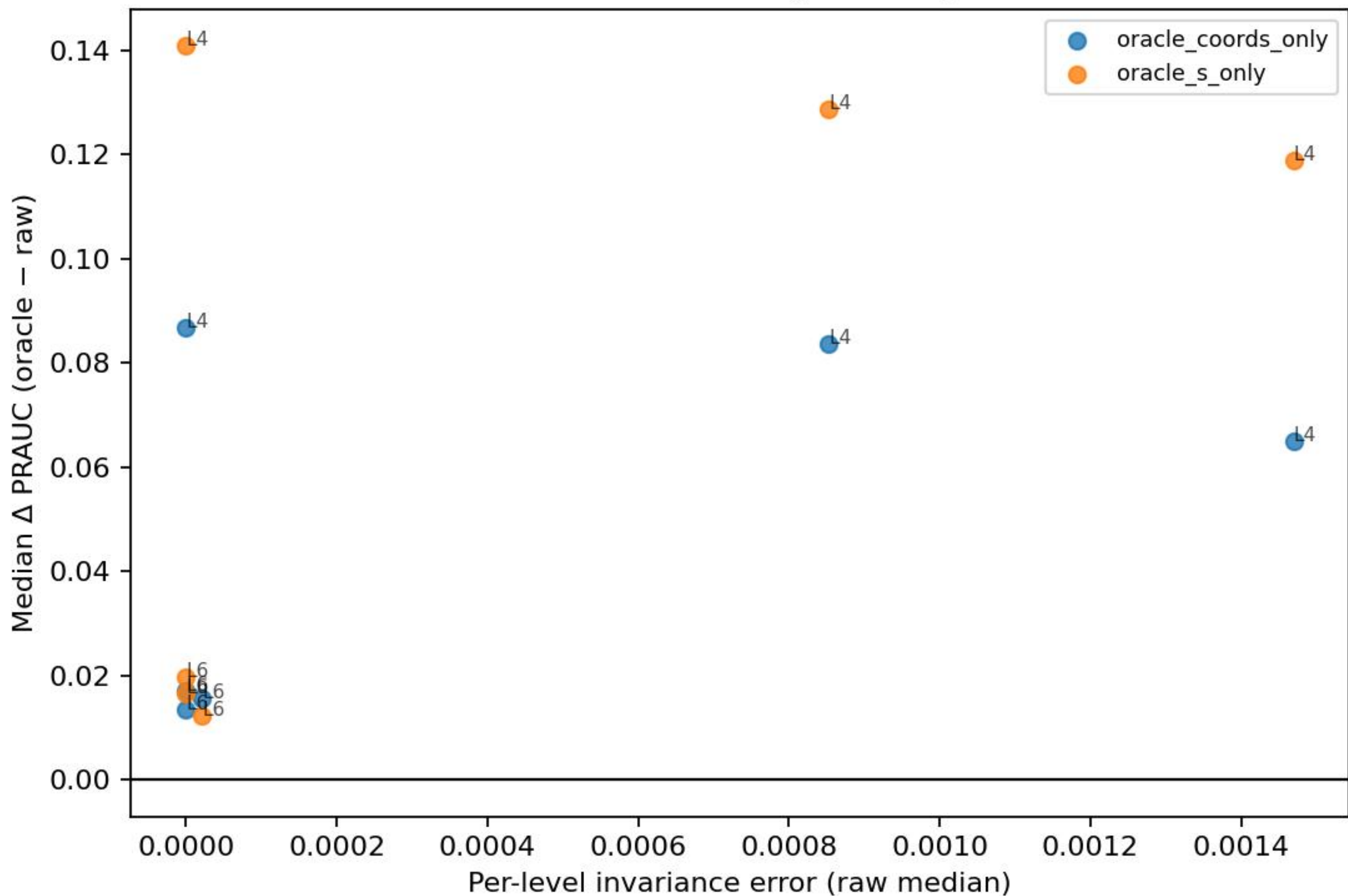Sum dominance vs generalization gap

No dominance metrics available

Invariance vs generalization gap

Δ vs invariance (per level)

No n-sweep data