

Text-Twin-Translation (T³): A Full-Stack Machine Learning Framework for Functional Material-Device Systems Discovery

Technical Appendices and Supplementary Material

A SUPPORTING FIGURES

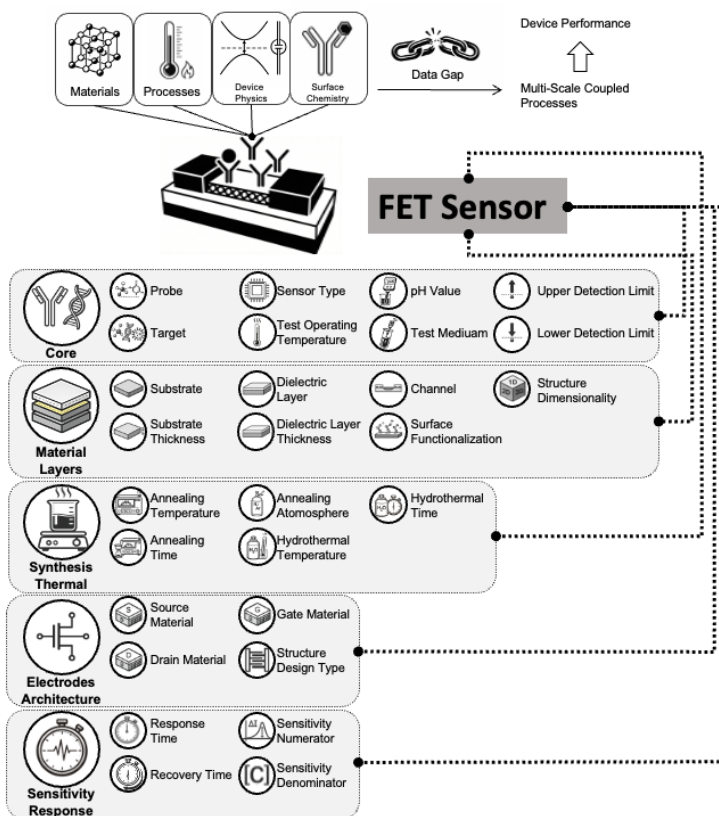


Figure S1: The conceptual gap between traditional “Fragmented View” (optimizing materials or physics in silos) and the reality of “Coupled Systems,” separated by the “Data Gap” of unstructured literature. This complexity is instantiated in FET sensors through a high-dimensional ontology mapping the intricate dependencies across synthesis, core chemistry, device architecture, and material layers that jointly determine performance.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

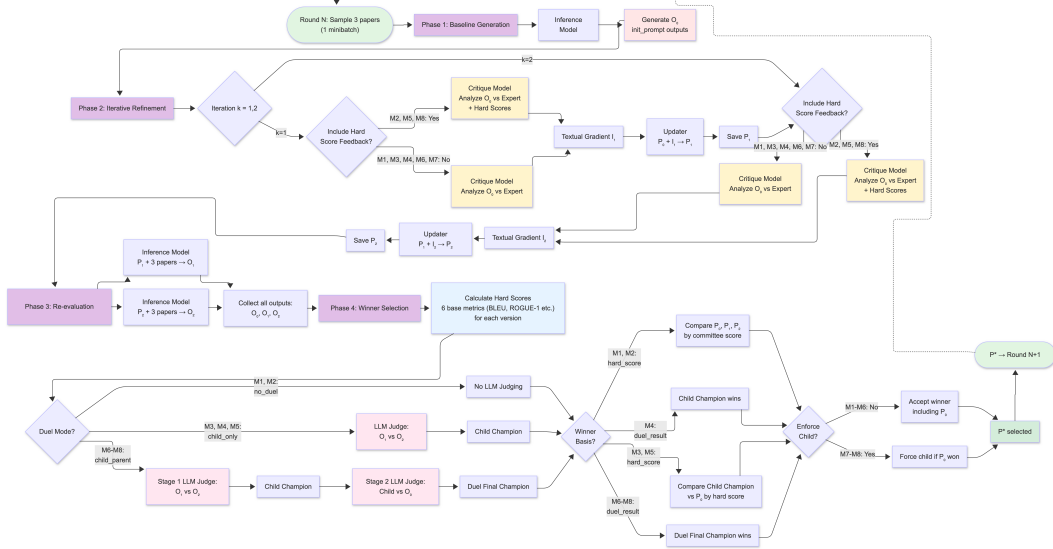


Figure S2: Flowchart schematic of the workflow of different modes: M1-M8 doing automatic prompt optimization.

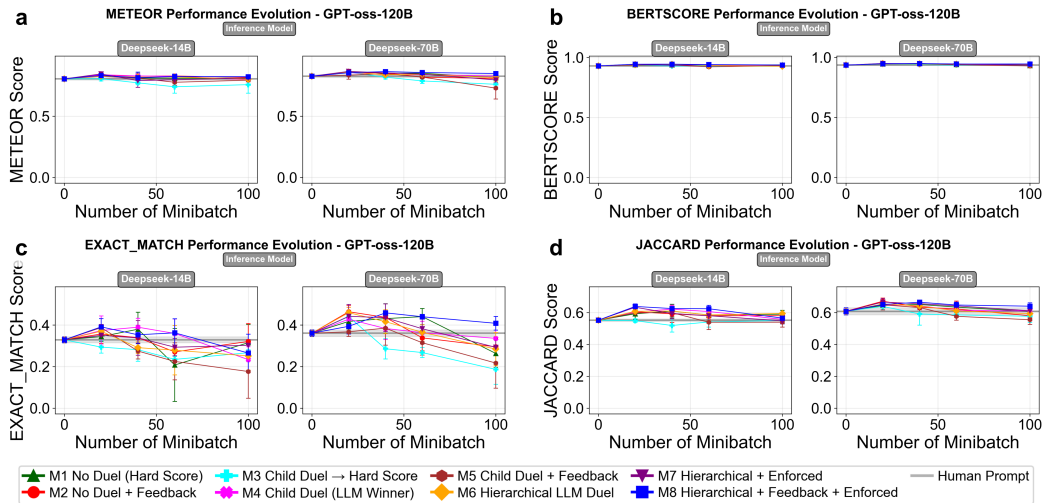


Figure S3: GPT-oss-120B evolution on four additional metrics (METEOR, BERTScore, Exact Match, Jaccard) averaged over 3 folds and inference models Deepseek-14B/70B.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

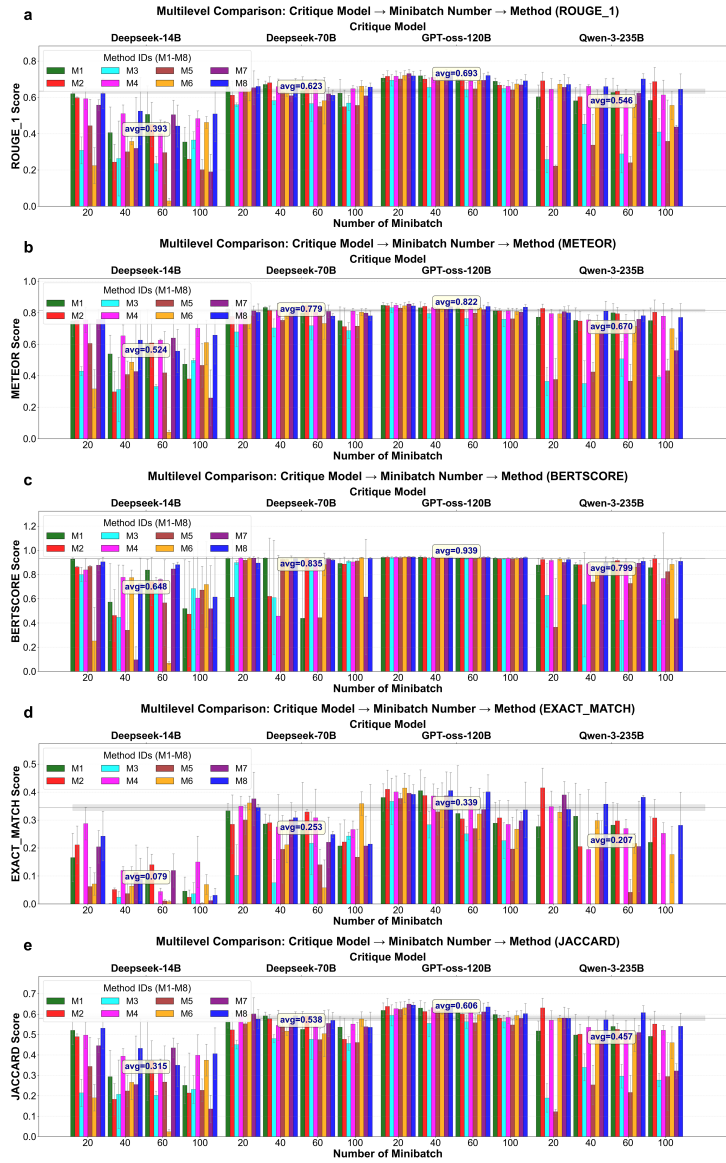
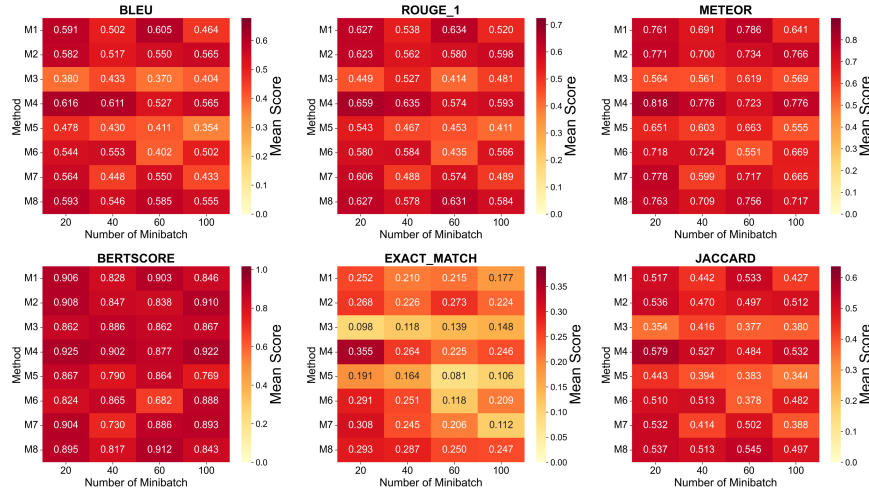


Figure S4: Multilevel barplots (critique → minibatch → method) for ROUGE-1, METEOR, BERTScore, Exact Match, and Jaccard across 3 folds (inference models averaged).

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

a

Performance Heatmaps for Inference Model: Deepseek-14B (Averaged over local critiques: Deepseek-14B, Deepseek-70B, Qwen-3-235B, GPT-oss-120B; 3 folds)



b

Performance Heatmaps for Inference Model: Deepseek-70B (Averaged over local critiques: Deepseek-14B, Deepseek-70B, Qwen-3-235B, GPT-oss-120B; 3 folds)

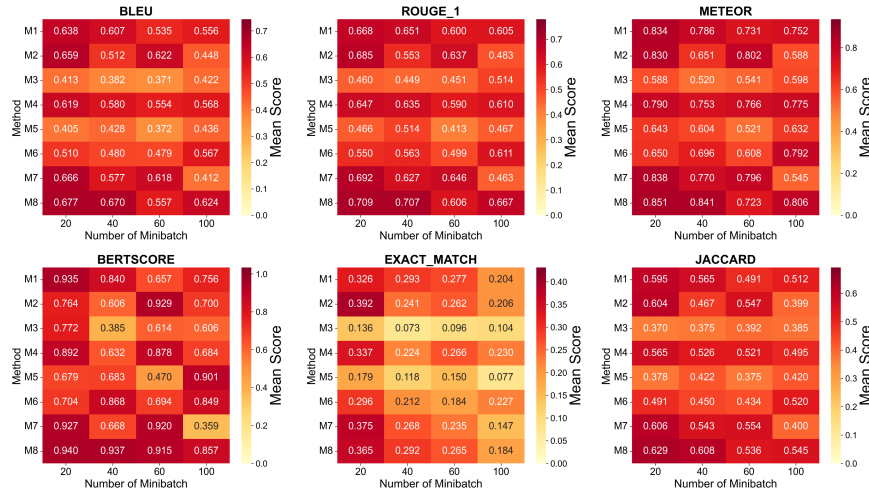


Figure S5: Per-inference-model heatmaps (Deepseek-14B and Deepseek-70B): methods M1–M8 × minibatch sizes (20/40/60/100) across six metrics, averaged over 3 folds.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

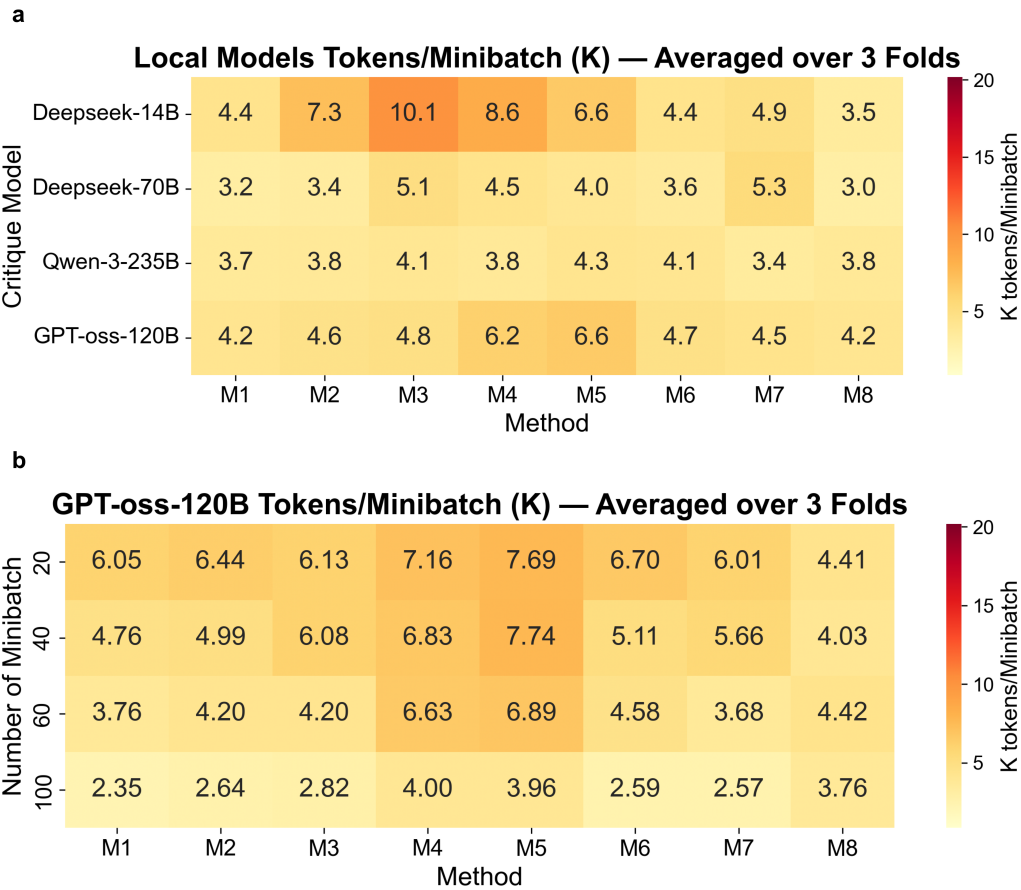


Figure S6: Token usage (K tokens per minibatch) for open-source critiques (top) and GPT-oss-120B (bottom), averaged over inference models and 3 folds.

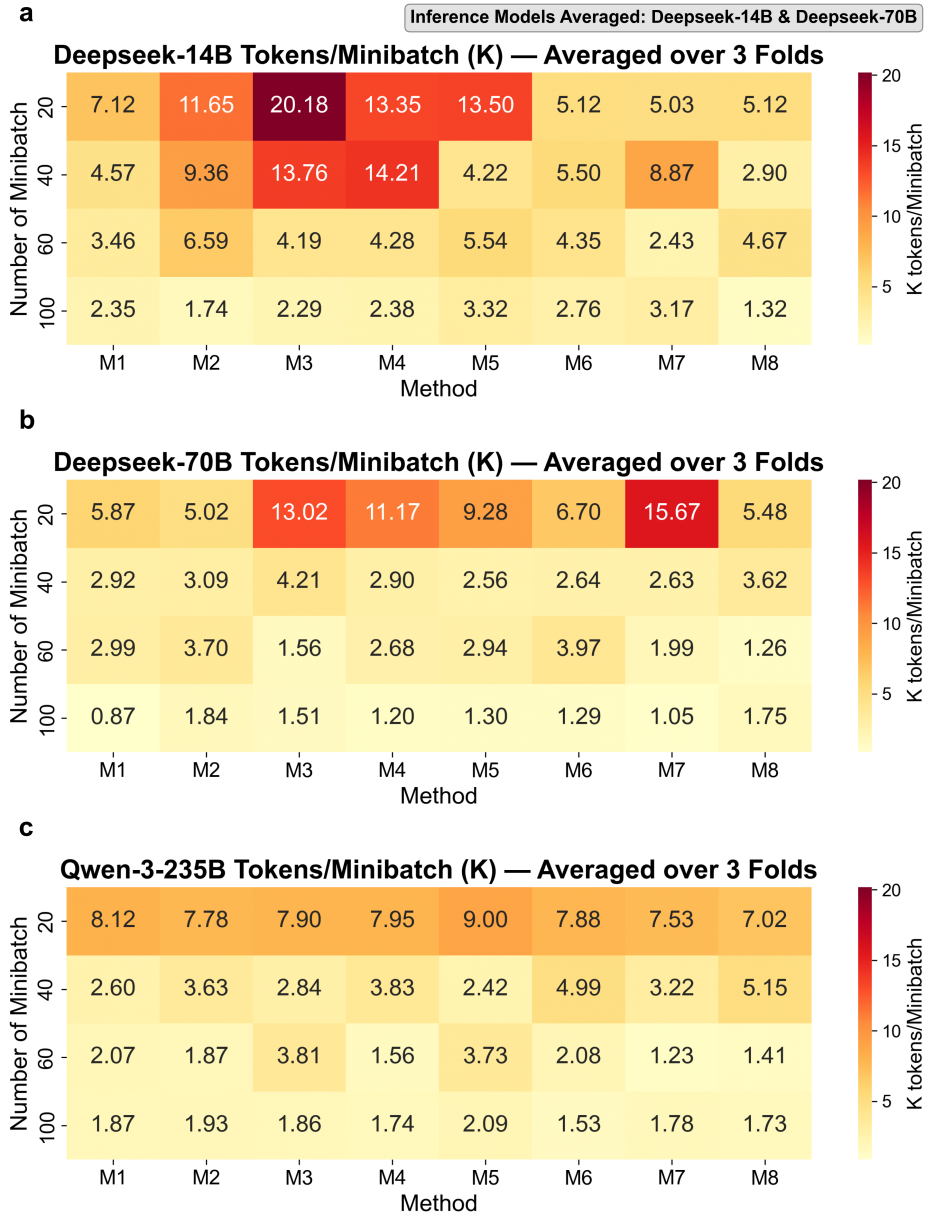


Figure S7: Token usage per minibatch (K tokens) for individual local critiques (Deepseek-14B/70B, Qwen-3-235B), averaged over inference models and 3 folds.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

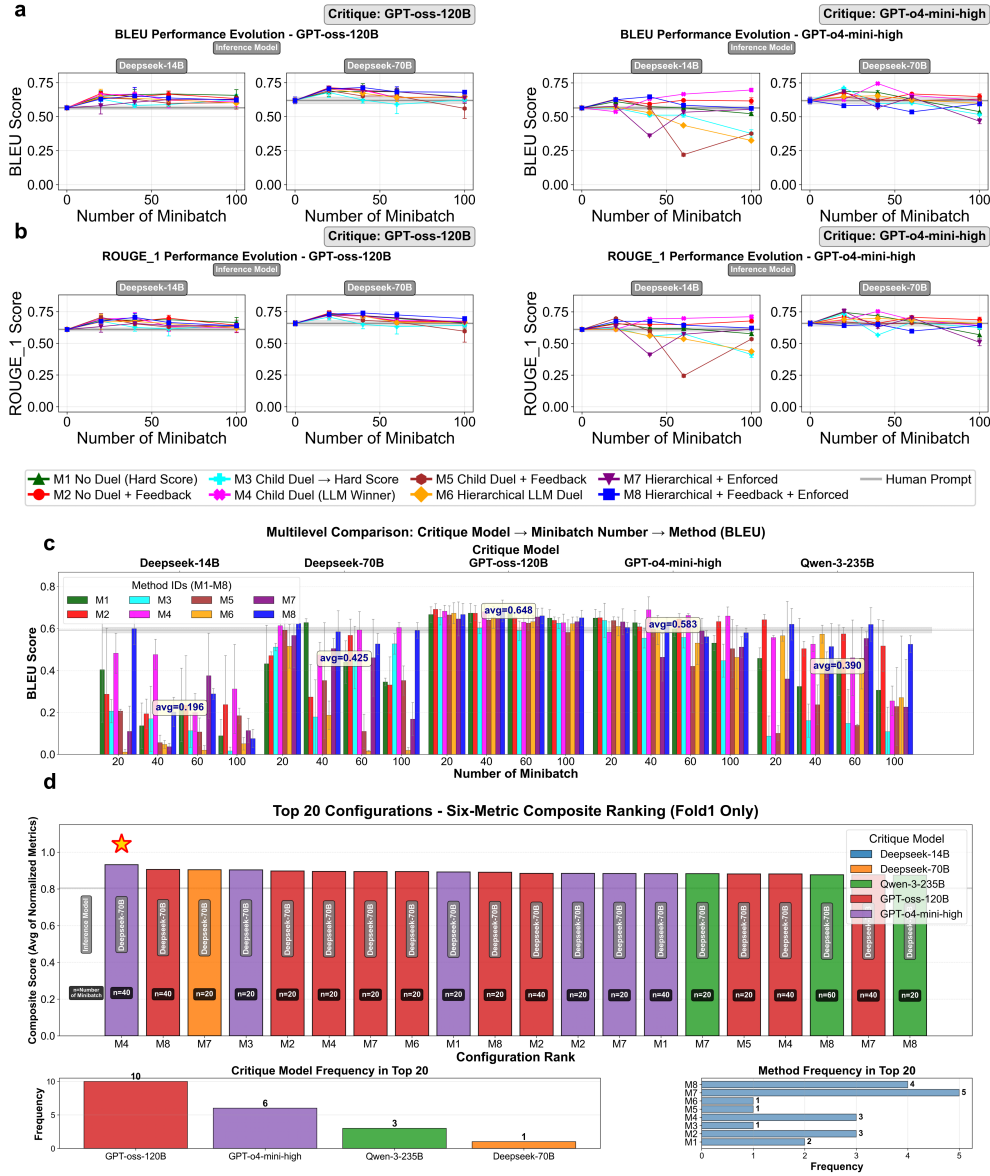


Figure S8: Fold-1 comparison of GPT-oss-120B vs. GPT-o4-mini-high: BLEU/ROUGE evolution, multilevel BLEU barplot, and final recommendation.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241

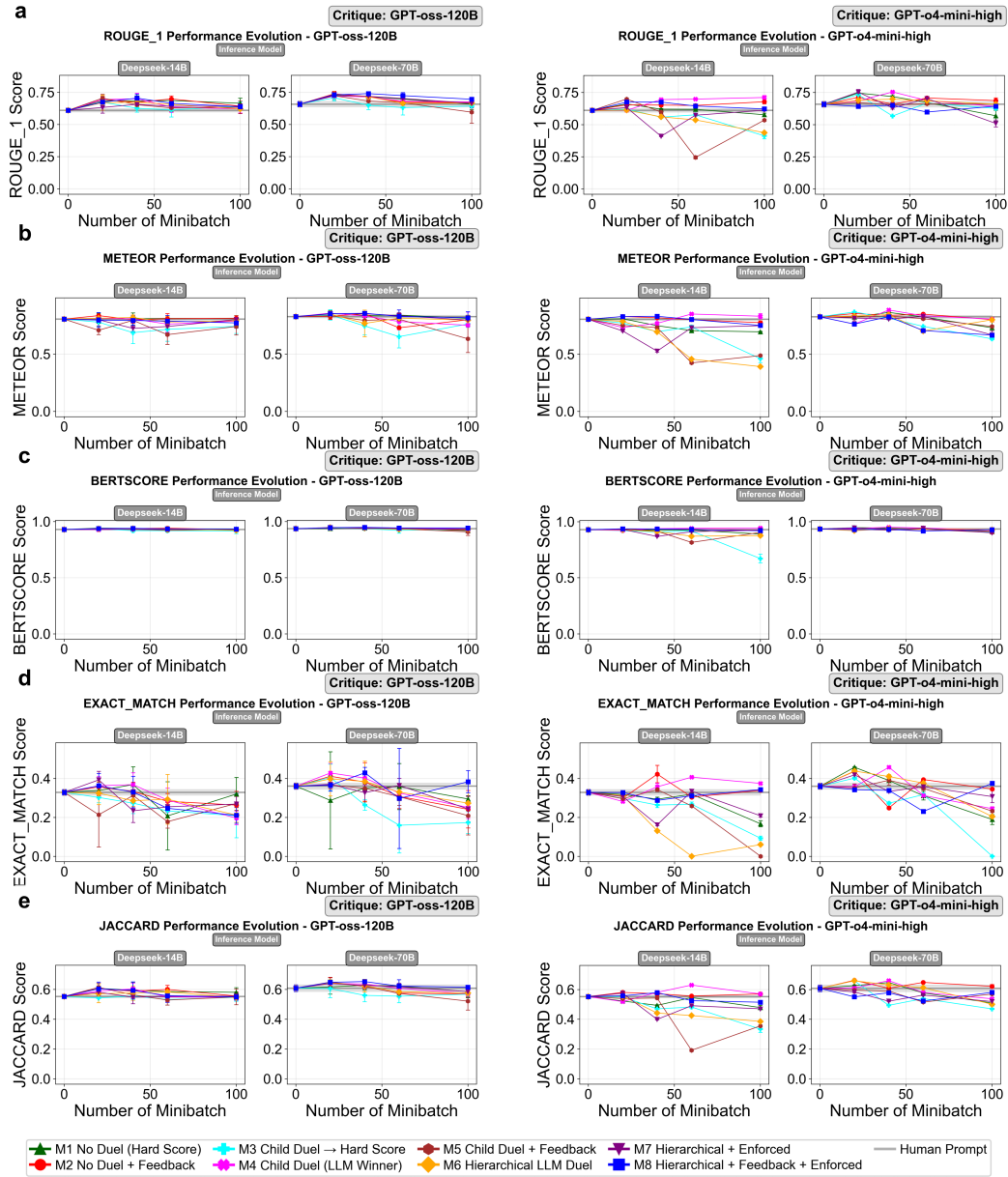


Figure S9: Fold-1 evolution comparisons for GPT-oss-120B vs. GPT-o4-mini-high on METEOR, BERTScore, Exact Match, and Jaccard.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

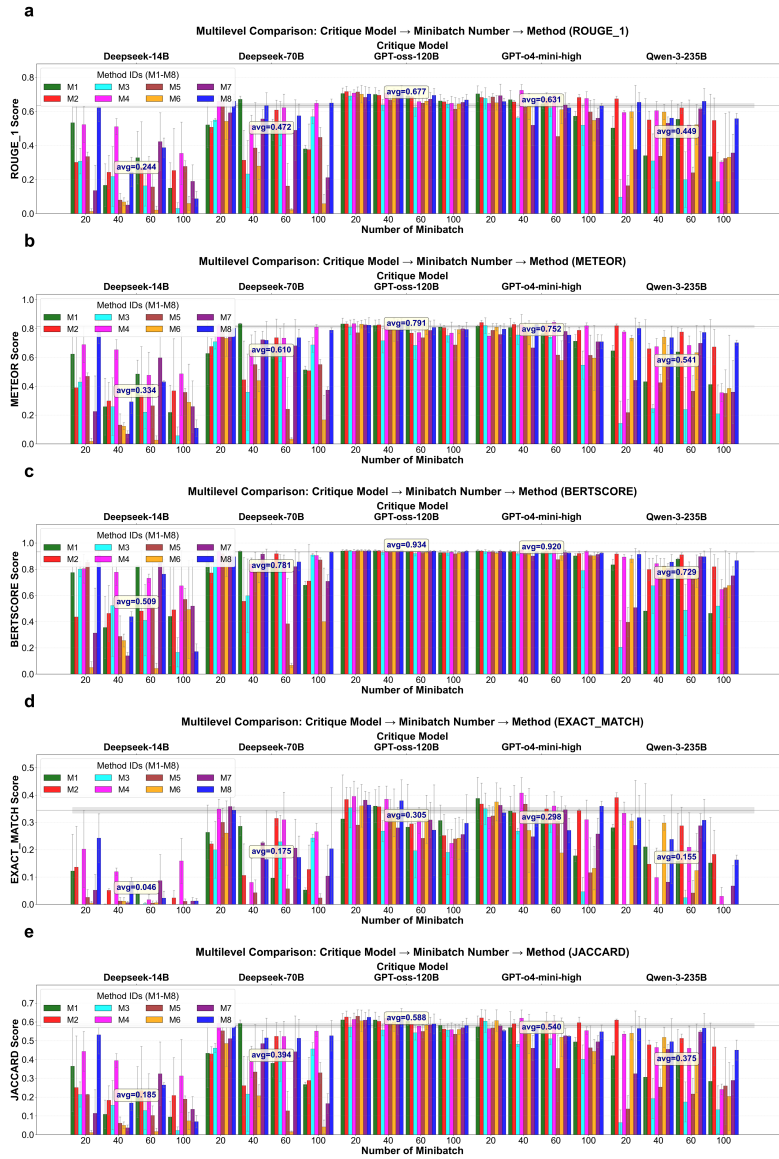


Figure S10: Fold-1 multilevel barplots (with GPT-o4-mini-high included) for ROUGE-1, METEOR, BERTScore, Exact Match, and Jaccard.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

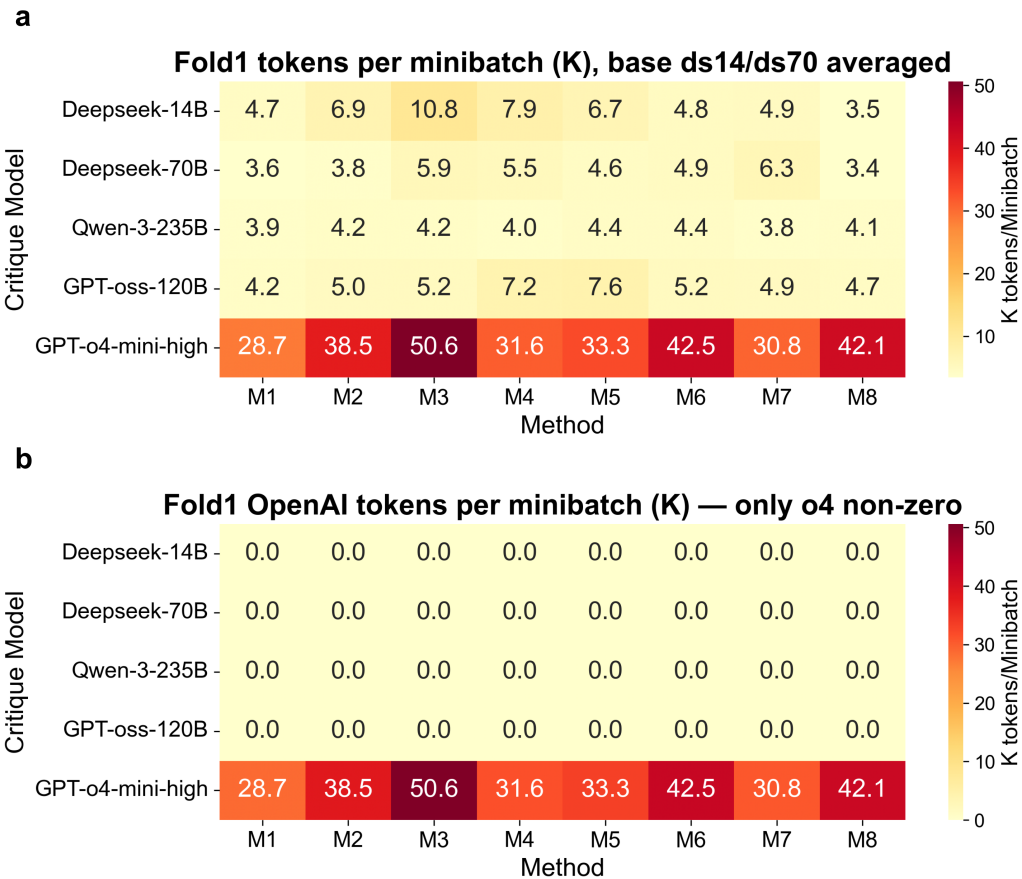


Figure S11: Fold-1 token cost heatmaps: total tokens per minibatch (left) and OpenAI-only tokens (right).

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

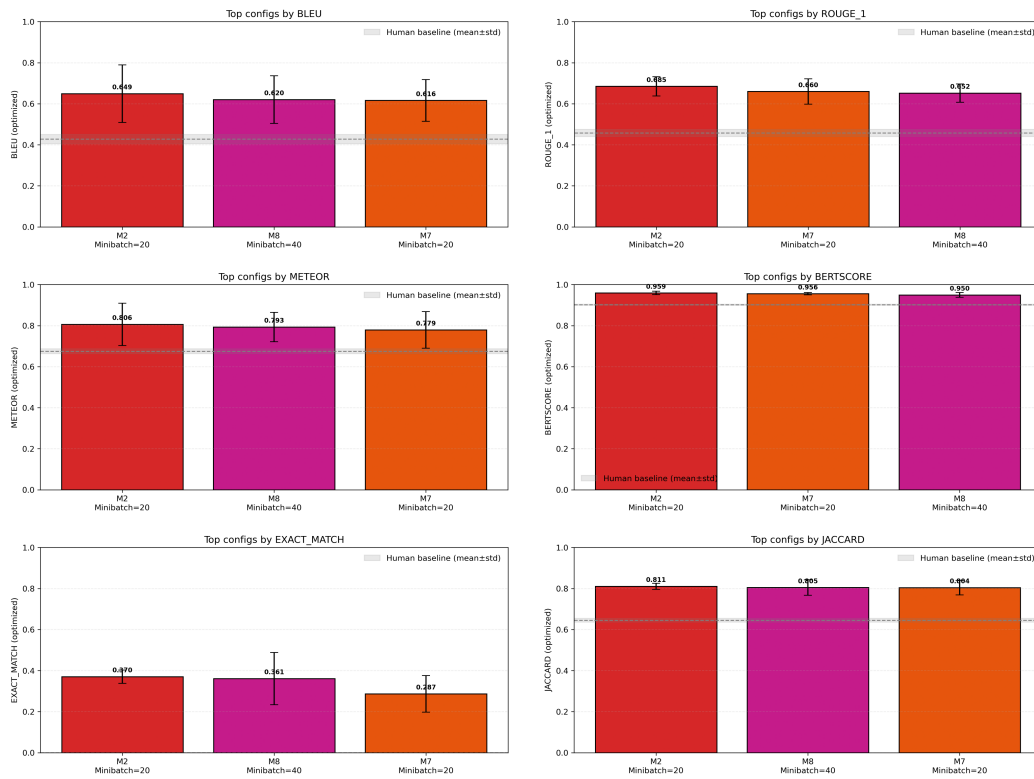


Figure S12: Extended field — Electrode Architecture: top configurations (M8-40, M7-20, M2-20) vs. human baseline across six metrics.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

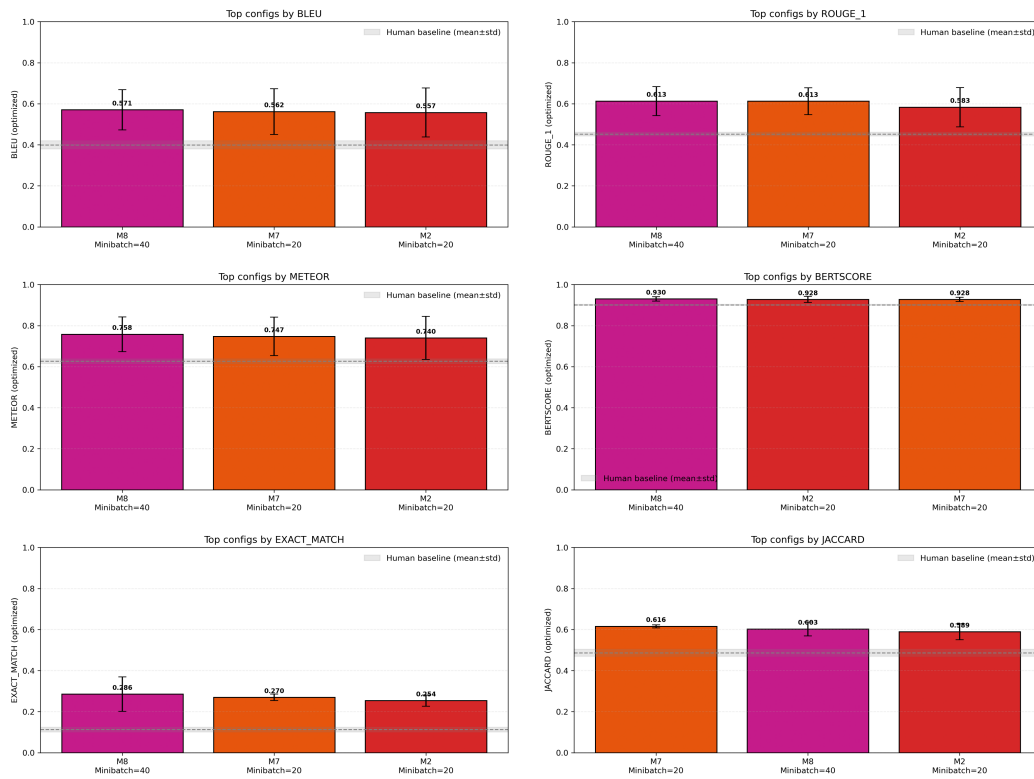


Figure S13: Extended field — Material Layer Composition: top configurations (M8-40, M7-20, M2-20) vs. human baseline across six metrics.

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

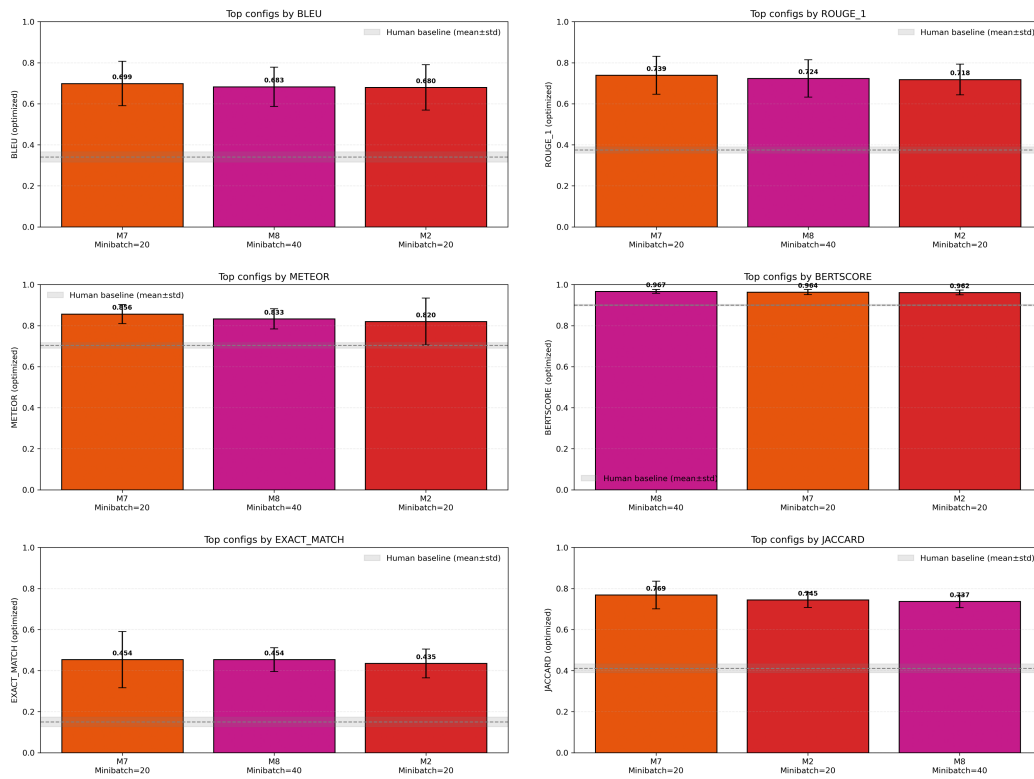


Figure S14: Extended field — Sensitivity and Response: top configurations (M8-40, M7-20, M2-20) vs. human baseline across six metrics.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

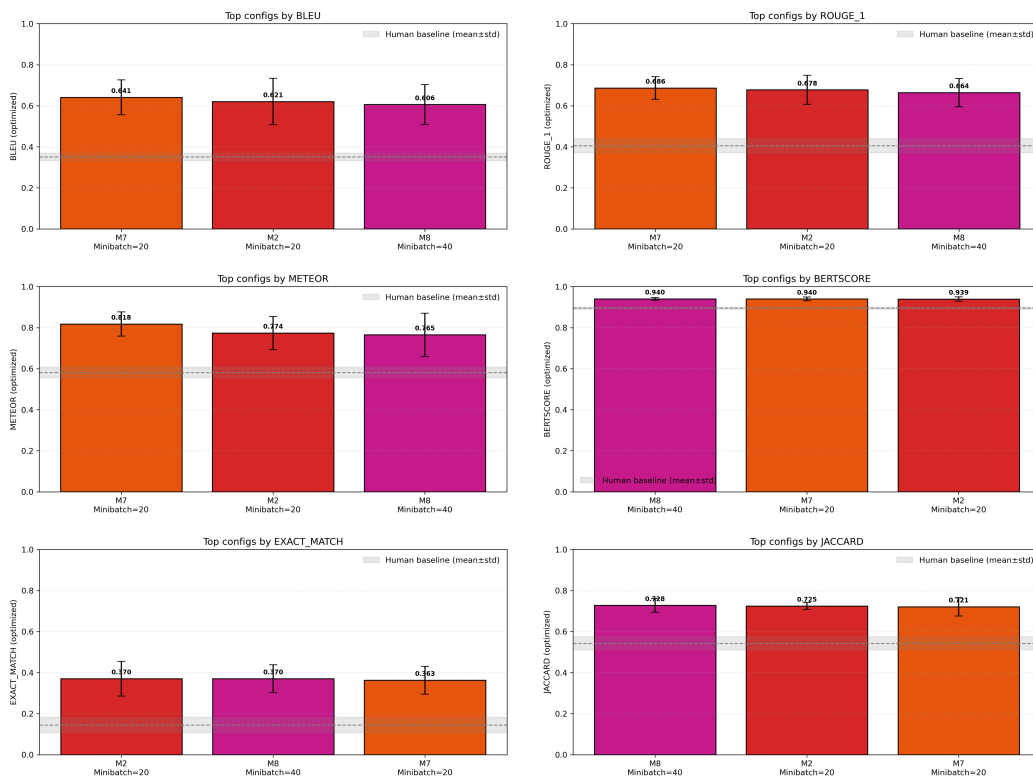


Figure S15: Extended field — Synthesis and Thermal Processing: top configurations (M8-40, M7-20, M2-20) vs. human baseline across six metrics.

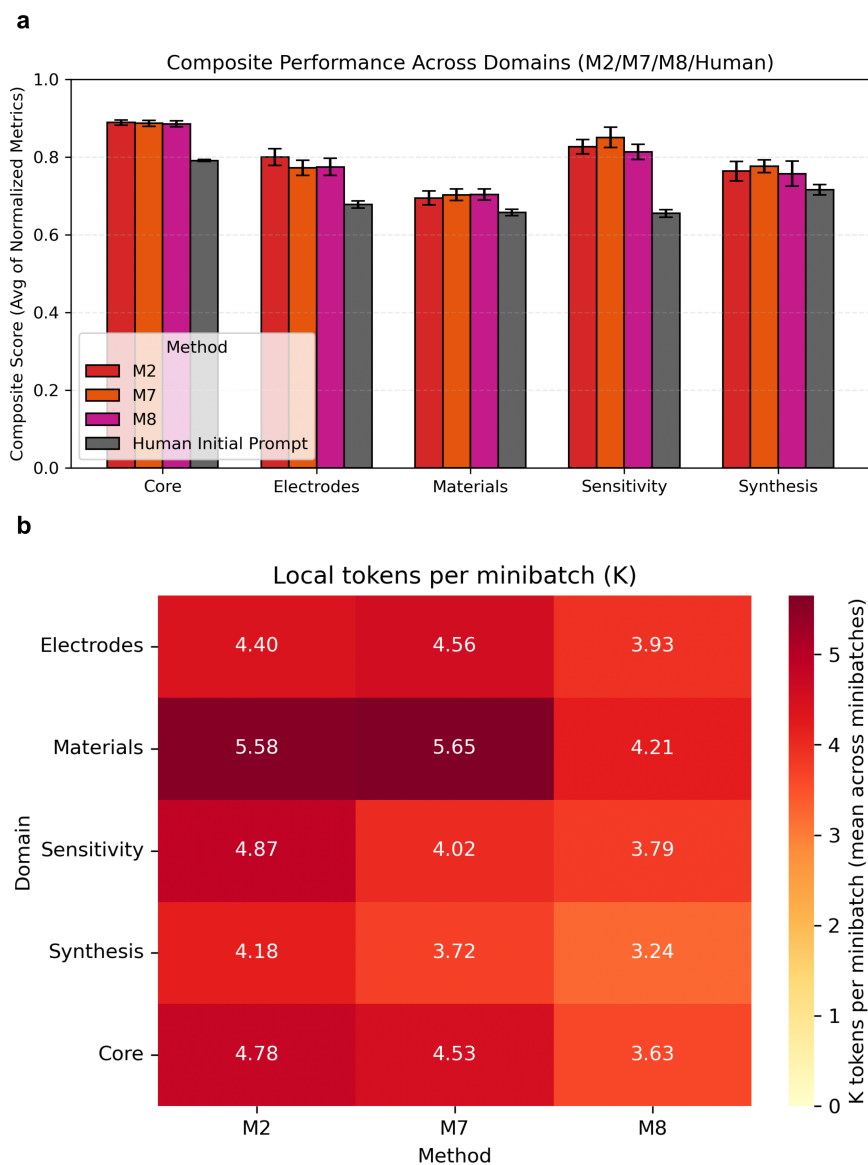


Figure S16: Extended fields summary (panels a–b). (a) Composite scores across five fields for the top-3 strategies vs. Human baseline. (b) Token usage per minibatch (K tokens) aggregated over the four extended fields for the top-3 strategies.

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

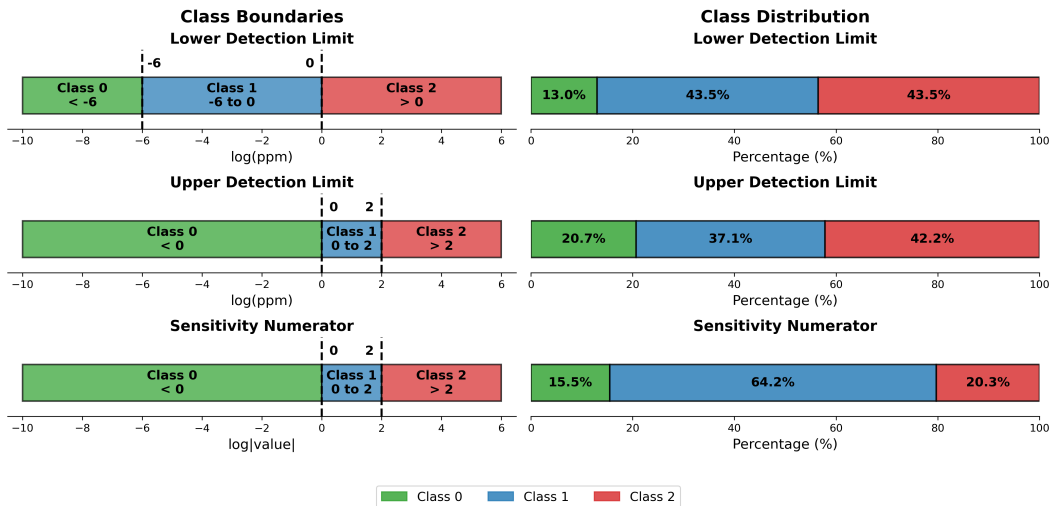


Figure S17: Classification task definitions for Twin-phase device performance prediction. Left: class boundaries on log scale for three prediction targets—Lower Detection Limit (LDL), Upper Detection Limit (UDL), and Sensitivity. Each target is discretized into three classes spanning multiple orders of magnitude. Right: class distribution percentages across the augmented dataset.

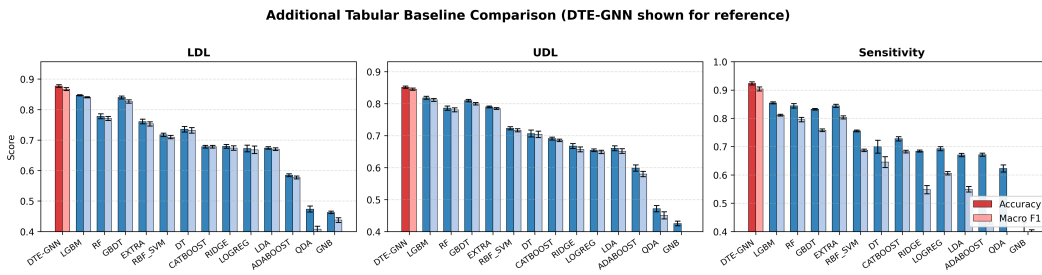


Figure S18: Additional tabular baseline comparison on held-out original data. Performance of 13 tabular machine learning models not shown in the main figure, across LDL, UDL, and Sensitivity prediction tasks. Models are ordered left-to-right by cross-task UCB ranking. While LightGBM (LGBM), Random Forest (RF), and Gradient Boosting (GBDT) achieve competitive performance, all tabular baselines underperform DTE-GNN (Figure 4a), confirming the advantage of heterogeneous graph representation.

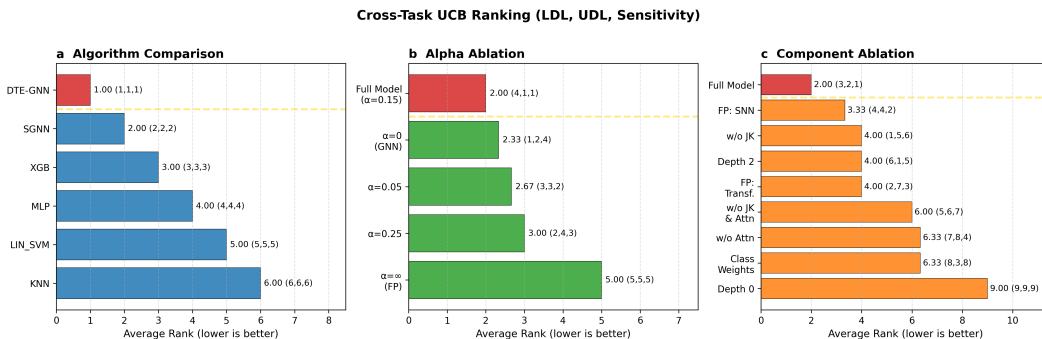


Figure S19: Cross-task UCB ranking visualization (see Section E for methodology). Horizontal bars show the average rank across LDL, UDL, and Sensitivity tasks for each configuration, with per-task ranks shown in parentheses. Lower average rank indicates more consistent top performance. (a) Algorithm comparison: DTE-GNN achieves perfect rank 1 across all tasks. (b) Alpha ablation: Full Model ($\alpha=0.15$) achieves the best average rank of 2.00. (c) Component ablation: Full Model achieves the best average rank of 2.00, with a 1.33 gap to the second-best configuration.

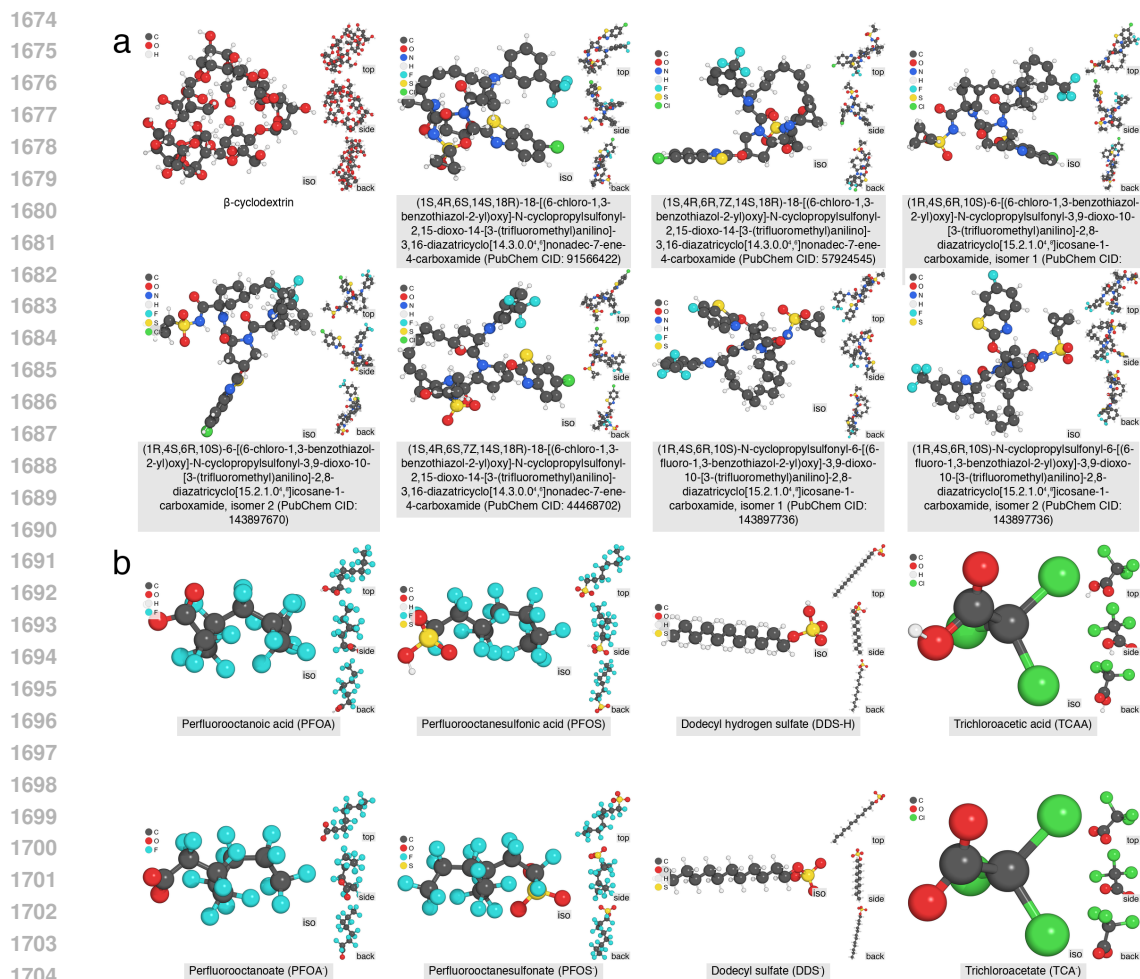
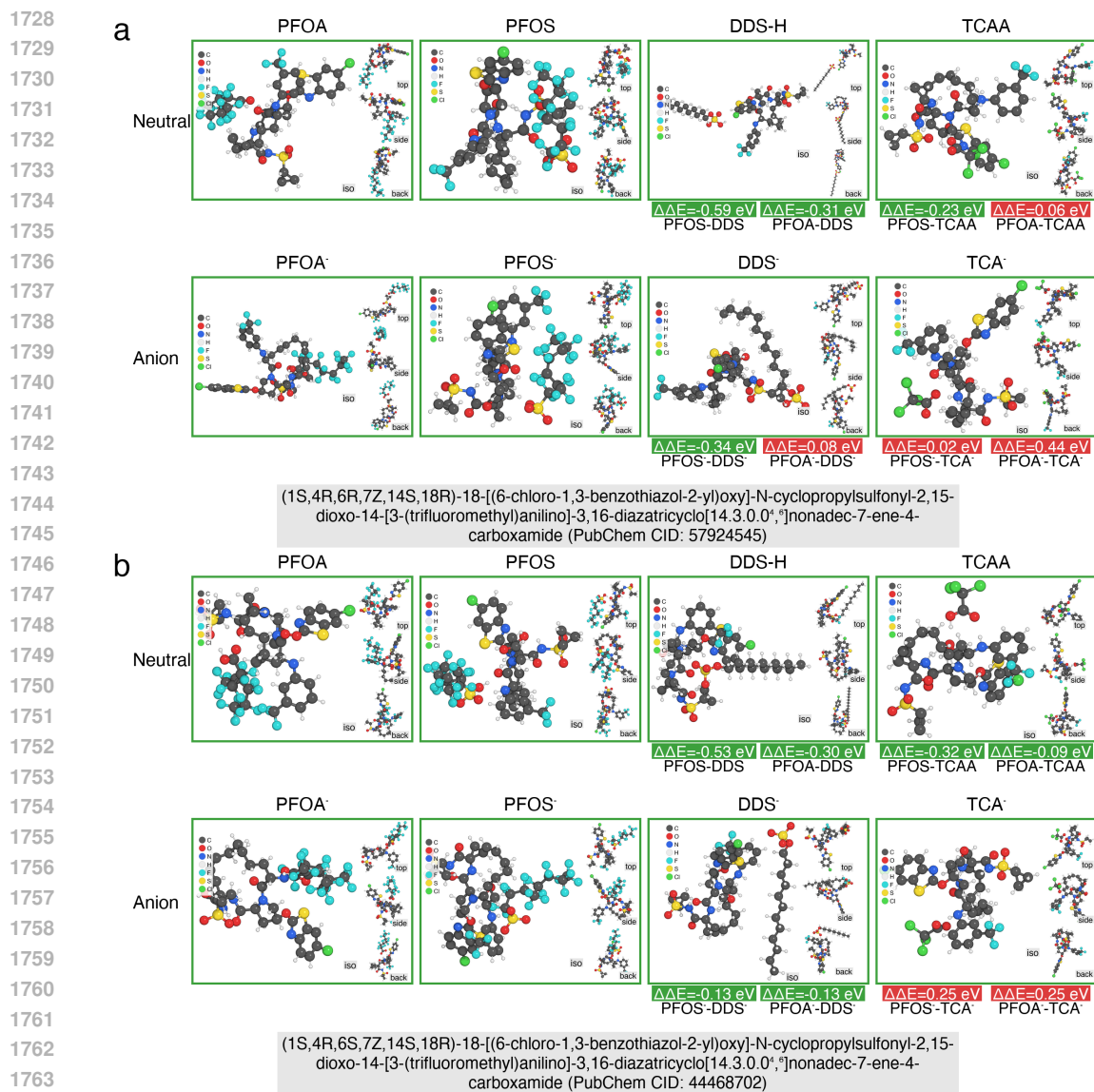
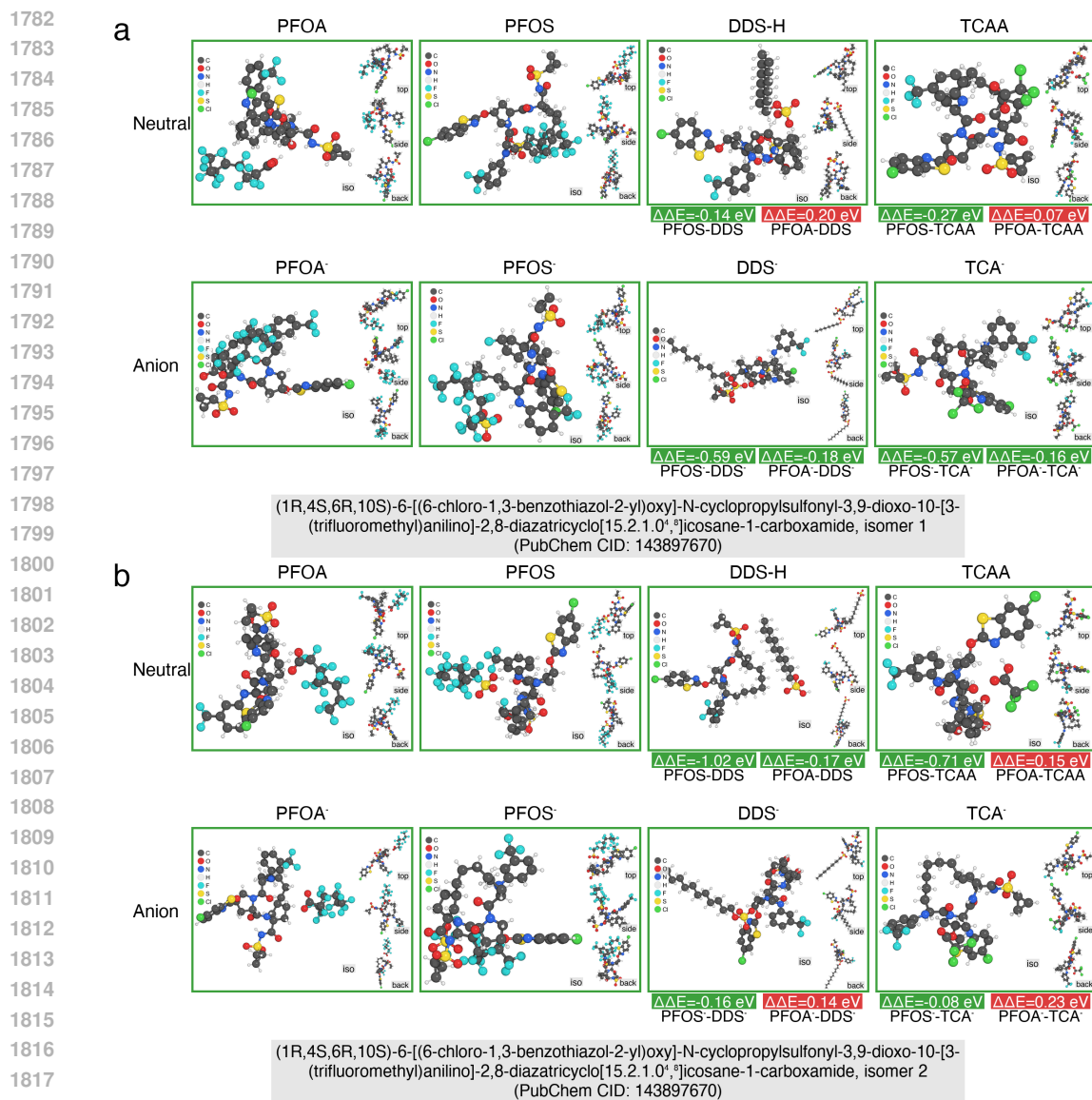


Figure S20: Structure views of the (a) probe and (b) target molecules/anions.



1765 Figure S21: Structure views and DFT-calculated binding energy difference (selectivity) of probe (a)
1766 CID-57924545 and (b) CID-44468702 over target molecules/anions.

1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781



1819 Figure S22: Structure views and DFT-calculated binding energy difference (selectivity) of probe (a)
1820 CID-143897670 isomer #1 and (b) CID-143897670 isomer #2 over target molecules/anions.

1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835

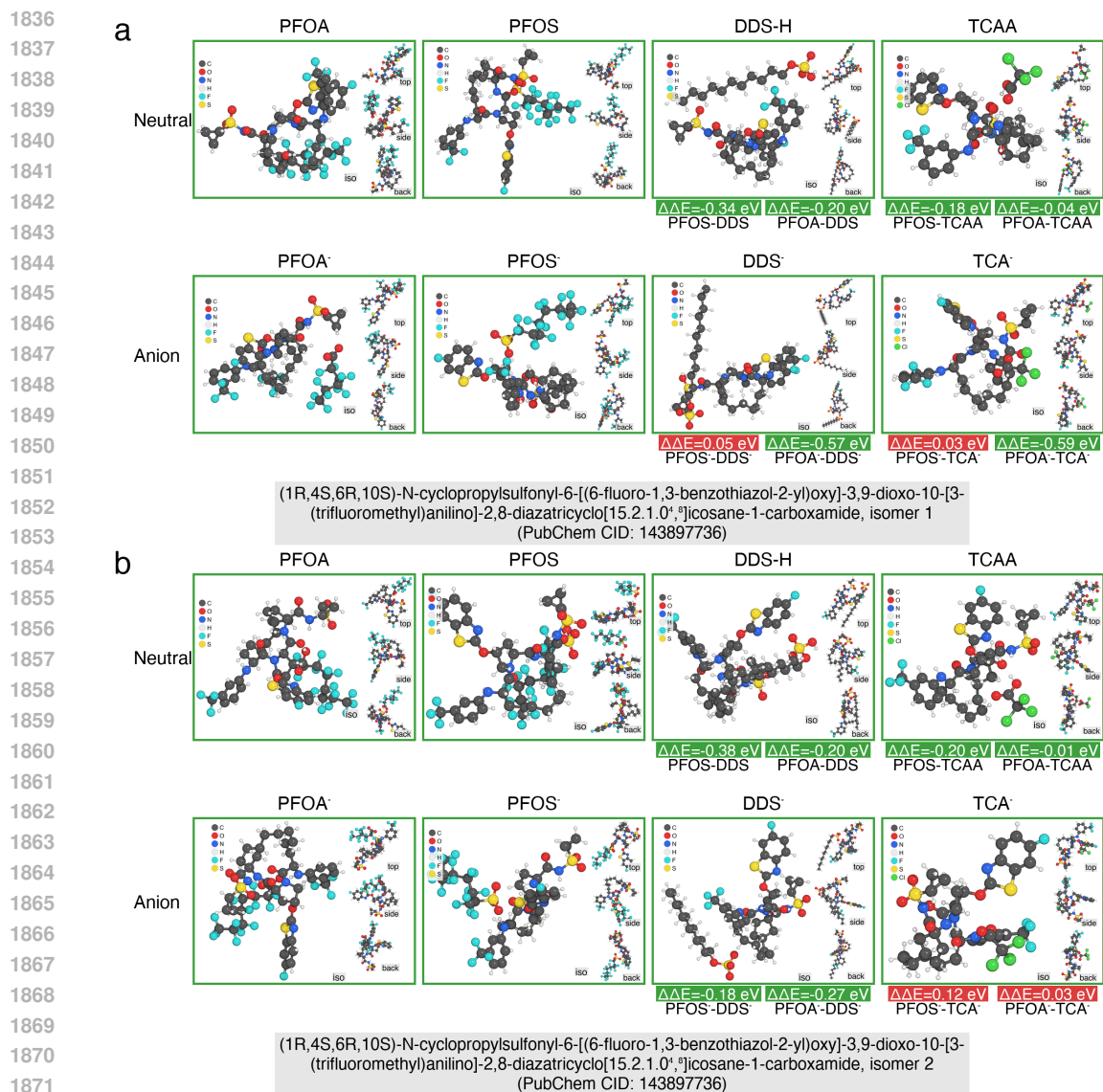


Figure S23: Structure views and DFT-calculated binding energy difference (selectivity) of probe (a) CID-143897736 isomer #1 and (b) CID-143897736 isomer #2 over target molecules/anions.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Table S1: DTE-GNN predicted device-performance scores and DFT binding-energy selectivity ($\Delta\Delta E$, eV) for the top five candidates and β -Cyclodextrin baseline. S_{target} (Eq. 8) is computed for PFOS. Negative $\Delta\Delta E$ indicates preferential PFAS binding (favorable selectivity). Three candidates (CIDs 44468702, 57924545, 91566422) are stereoisomers with identical 2D fingerprints and indistinguishable DTE-GNN scores. Two candidates (CIDs 143897670 and 143897736) have two isomers. Bold marks the only candidate with consistently negative $\Delta\Delta E$ across all conditions.

Probe	DTE-GNN (PFOS)				$\Delta\Delta E$ Neutral (eV)				$\Delta\Delta E$ Anion (eV)			
	P_0^{LDL}	P_2^{LDL}	P_2^{Sens}	S_{target}	PFOS-DDS	PFOA-DDS	PFOS-TCAA	PFOA-TCAA	PFOS ⁻ -DDS ⁻	PFOA ⁻ -DDS ⁻	PFOS ⁻ -TCA ⁻	PFOA ⁻ -TCA ⁻
β -Cyclodextrin	0.998	0.000	0.003	~ 0	-1.90	-0.43	-1.57	-0.10	+0.54	-0.21	+0.68	-0.07
CID-143897670-isomer1	0.912	0.653	0.200	0.130	-0.14	+0.20	-0.27	+0.07	-0.59	-0.18	-0.57	-0.16
CID-143897670-isomer2	0.912	0.653	0.200	0.130	-1.02	-0.17	-0.71	+0.15	-0.16	+0.14	-0.08	+0.23
CID-143897736-isomer1	0.918	0.585	0.193	0.113	-0.34	-0.20	-0.18	-0.04	+0.05	-0.57	+0.03	-0.59
CID-143897736-isomer2	0.918	0.585	0.193	0.113	-0.38	-0.20	-0.20	-0.01	-0.18	-0.27	+0.12	+0.03
CID-91566422	0.895	0.783	0.138	0.105	-1.53	-1.04	-1.12	-0.62	-0.31	-0.53	-0.23	-0.46
CID-44468702	0.895	0.783	0.138	0.105	-0.53	-0.30	-0.32	-0.09	-0.13	-0.13	+0.25	+0.25
CID-57924545	0.895	0.783	0.138	0.105	-0.59	-0.31	-0.23	+0.06	-0.34	+0.08	+0.02	+0.44

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

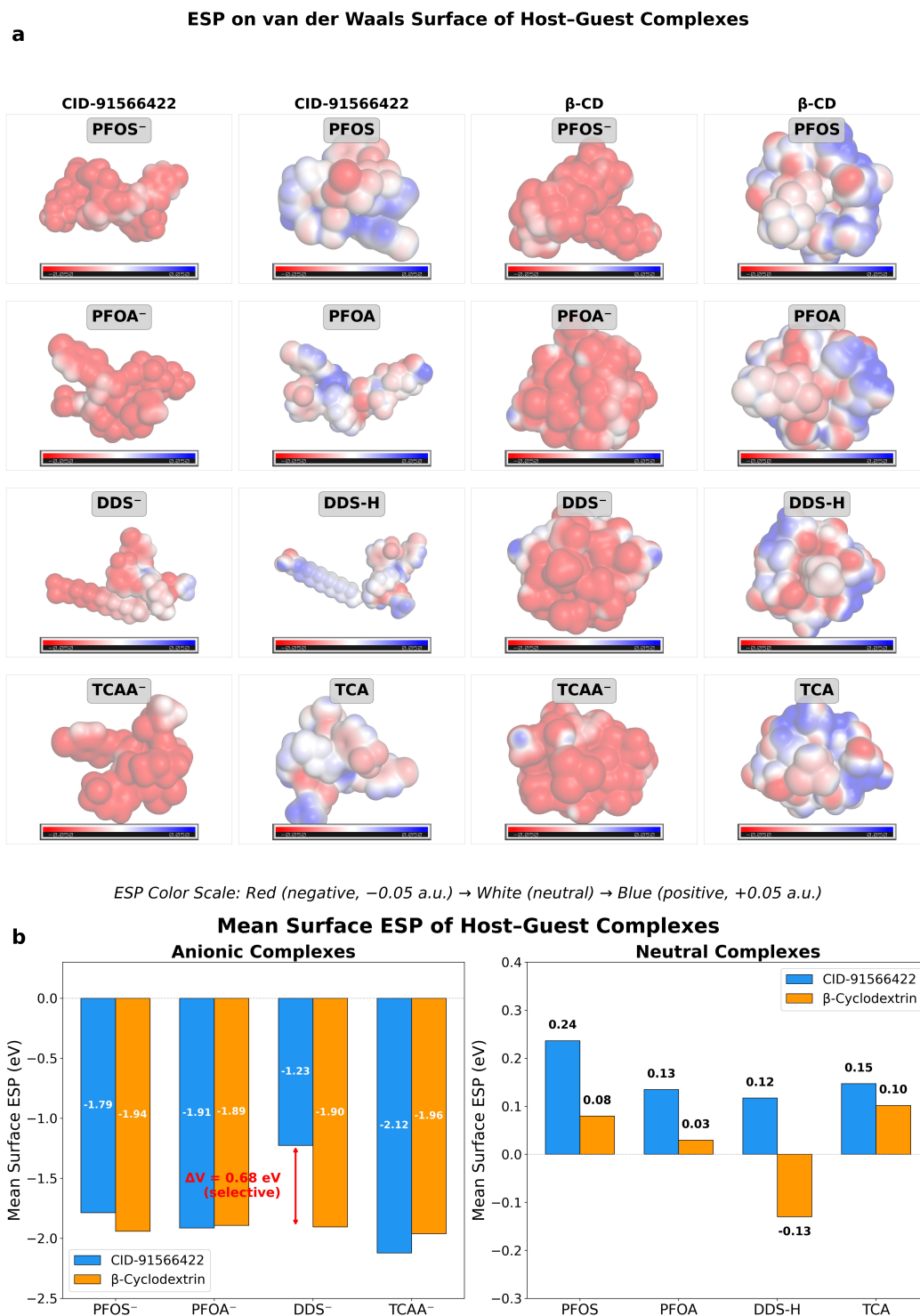


Figure S24: **Electrostatic potential (ESP) analysis of host-guest complexes.** (a) ESP mapped onto the van der Waals surface ($\rho = 0.001$ a.u.) for CID-91566422 and β -Cyclodextrin complexes with four guests under anionic and neutral conditions. Color scale: red (-0.05 a.u.) \rightarrow white (0) \rightarrow blue ($+0.05$ a.u.). (b) Mean surface ESP (\bar{V} , eV) comparison. CID-91566422 shows a distinctly less negative \bar{V} for DDS⁻ (-1.23 eV) compared to PFAS guests (-1.79 to -2.12 eV), while β -Cyclodextrin shows uniform $\bar{V} \approx -1.9$ eV for all guests.

B ADDITIONAL RELATED WORK

LLM-BASED SCIENTIFIC INFORMATION EXTRACTION

Scientific information extraction (IE) has long sought to convert the narrative style of papers into machine-actionable records. Earlier materials/chemistry pipelines relied on rule-based parsing and domain heuristics, exemplified by ChemDataExtractor [Swain & Cole \(2016\)](#). With the rise of scientific-language pretraining, transformer encoders such as SciBERT [Beltagy et al. \(2019\)](#) enabled stronger scientific NER and sentence-level semantics, while domain-focused IE systems further emphasized entity normalization and scalable extraction in materials literature [Weston et al. \(2019\)](#). Beyond NER, end-to-end scientific IE has also been framed as joint extraction of entities, relations, and coreference to support structured scholarly representations [Luan et al. \(2018\)](#). Despite these advances, most approaches either require substantial labeled data or must be carefully adapted to each narrow schema and subdomain, a challenge that becomes acute for device papers where fabrication, geometry, interfaces, and performance metrics are tightly entangled.

More recently, generative LLMs have been used for schema-driven structured extraction from scientific text, including large-scale demonstrations in materials and chemistry [Dagdelen et al. \(2024\)](#). However, these systems often depend on extensive fine-tuning, brittle prompt crafting, or post-hoc repair to control output validity and cross-field consistency. Domain-specific pretraining can help but does not eliminate the supervision bottleneck for specialized schemas [Trewartha et al. \(2022\)](#). In contrast, our T^3 framework is designed as an end-to-end pipeline that couples autonomous, prompt-optimized extraction with downstream device-level prediction, directly targeting the extraction-modeling gap for heterogeneous FET sensor literature.

PROMPT OPTIMIZATION AND TEXTUAL GRADIENTS

Prompting has become a central mechanism for steering frozen LLMs, including reasoning-oriented prompting strategies such as chain-of-thought [Wei et al. \(2022\)](#). To reduce reliance on manual prompt engineering, automated prompt search and refinement has been explored via LLM-driven proposal-and-selection loops. Automatic Prompt Engineer (APE) [Zhou et al. \(2023\)](#) and Optimization by PROMpting (OPRO) [Yang et al. \(2024\)](#) treat prompts as optimizable parameters, iteratively generating candidate instructions and selecting those that improve task performance. In parallel, system-building frameworks such as DSPy [Khattab et al. \(2023\)](#) compile multi-step LLM pipelines and tune prompts programmatically, highlighting a shift from single-prompt tuning to compositional optimization.

A complementary line of work introduces optimization signals that resemble gradients, where critiques or losses are converted into text edits. Automatic Prompt Optimization (APO) operationalizes this idea through minibatch "natural language gradients" plus guided search [Pryzant et al. \(2023\)](#), while TextGrad generalizes the concept by backpropagating LLM-generated feedback through computation graphs to improve upstream components [Yuksekgonul et al. \(2025b\)](#). Most evaluations, however, emphasize generic NLP/reasoning benchmarks, leaving open how to robustly optimize prompts under scientific constraints (units, normalization, nested schemas, and strict output validity). Our contribution extends textual-gradient optimization into a multi-mode TextGrad regime tailored to scientific IE, enabling reliable schema-conformant extraction with minimal human annotation and directly powering the Text \rightarrow Twin translation in T^3 .

GRAPH NEURAL NETWORKS FOR MATERIALS AND MOLECULES

Graph neural networks (GNNs) are now the dominant paradigm for molecular and materials property prediction because they encode relational inductive biases over atoms and bonds. Message Passing Neural Networks (MPNNs) provide a unifying formulation [Gilmer et al. \(2017\)](#), while geometric architectures such as SchNet [Schütt et al. \(2017\)](#) and directional message passing (DimeNet) [Klicpera et al. \(2020\)](#) capture distance- and angle-dependent interactions critical for quantum properties. In materials science, crystal graph models extend message passing to periodic solids [Xie & Grossman \(2018\)](#), and more recent variants incorporate richer geometric primitives (e.g., line graphs) to improve accuracy on benchmark property prediction tasks [Choudhary & DeCost \(2021\)](#). Alongside GNNs, learned and hand-crafted fingerprints remain widely used across modalities, including

2052 extended-connectivity fingerprints for molecules [Rogers & Hahn \(2010\)](#) and composition-based fea-
2053 turizations for inorganic materials [Ward et al. \(2016\)](#).

2054
2055 Despite strong performance at the molecule/crystal level, most GNN work assumes a single con-
2056 nected structure as the prediction object, whereas real sensing devices are multi-component sys-
2057 tems whose performance emerges from topology, interfaces, and coupled functional modules.
2058 Generic relational modeling (e.g., multi-relation GNNs) can represent typed edges and interactions
2059 [Schlichtkrull et al. \(2018b\)](#), but has rarely been specialized to device-scale, physics-informed topolo-
2060 gies. Our Twin component closes this gap by constructing a physics-informed heterogeneous graph
2061 for FET sensors (device topology + cross-domain fingerprints), enabling device-level predictions
2062 rather than single-material property inference.

2063 2064 2065 FET SENSOR MODELING AND PERFORMANCE PREDICTION

2066
2067
2068 FET sensors and related BioFET/ISFET devices have traditionally been analyzed using physics-
2069 based models that connect surface charge, electrostatics, and carrier transport to measurable signals.
2070 Foundational ISFET work established the core sensing principle in transistor form [Bergveld \(1970\)](#),
2071 and subsequent analyses studied fundamental limits arising from screening and electrolyte effects
2072 [Stern et al. \(2007\)](#) as well as performance ceilings for nanoscale BioFET sensors [Nair & Alam](#)
2073 [\(2006\)](#). These approaches are physically interpretable, but high-fidelity simulation and parameter
2074 calibration can be costly and difficult to scale across diverse materials, geometries, and function-
2075 alization strategies; moreover, device physics references emphasize that performance depends on
2076 coupled choices across stacks and interfaces rather than isolated parameters [Sze & Ng \(2006\)](#).

2077
2078 In response, data-driven methods have been increasingly used for sensor modeling and inference,
2079 including ML for biosensor signal processing and performance characterization [Cui et al. \(2020\)](#).
2080 Yet most learning-based pipelines flatten devices into tabular features, discarding explicit topology
2081 and making it hard to generalize across architectures (e.g., gate stacks, channel classes, probe layers)
2082 or to propagate fabrication/process effects through coupled interfaces. Building on prior graph-based
2083 progress for FET chemical sensor prediction [Ferreira et al. \(2025\)](#), T³ introduces a device-topology-
2084 aware, physics-informed heterogeneous representation that can be populated automatically from
2085 literature and trained end-to-end for device-level performance prediction.

2086 2087 SCIENTIFIC KNOWLEDGE GRAPHS AND LITERATURE MINING

2088
2089
2090 Scientific knowledge graphs (KGs) and scholarly corpora aim to represent literature at scale to sup-
2091 port retrieval, reasoning, and discovery. Large infrastructure efforts have produced heterogeneous
2092 literature graphs [Ammar et al. \(2018\)](#) and open corpora with structured full text and metadata [Lo](#)
2093 [et al. \(2020\)](#), providing foundations for downstream IE and linking. In materials science, major open
2094 databases such as the Materials Project [Jain et al. \(2013\)](#) and AFLOW [Curtarolo et al. \(2012\)](#) have
2095 accelerated property-driven research, but they primarily capture computed or curated structured data
2096 rather than the full experimental detail embedded in narrative papers.

2097
2098 To complement curated resources, literature mining has generated domain-specific datasets by
2099 extracting experimental procedures and properties directly from text, including synthesis recipe
2100 databases [Kononova et al. \(2019\)](#) and auto-generated magnetic transition temperature resources
2101 [Court & Cole \(2018\)](#). Beyond symbolic KGs, representation learning over text (e.g., word em-
2102 beddings) has been shown to recover latent materials knowledge and even anticipate discoveries
2103 [Tshitoyan et al. \(2019\)](#), and surveys have emphasized the broader opportunity and challenges of
2104 NLP/IE for materials databases [Olivetti et al. \(2020\)](#). A persistent limitation, however, is that ex-
2105 tracted entities and relations are often not wired into downstream predictive models—KGs remain
largely descriptive. T³ directly addresses this disconnect by designing extraction schemas that in-
stantiate a predictive device-twin graph, linking literature mining to topology-aware performance
models in a single end-to-end pipeline.

C METHOD AND TECHNICAL DETAILS

C.1 METHOD AND TECHNICAL DETAILS FOR “TEXT”

C.1.1 DATASET AND TASK DEFINITION

FET Sensor Literature Corpus Building upon prior work on neuromorphic spiking graph neural networks for FET sensor design [Ferreira et al. \(2025\)](#), we leverage a curated collection of 844 peer-reviewed scientific publications on field-effect transistor (FET) sensors. These articles were sourced from major scientific databases including IEEE Xplore, ACS Publications, Royal Society of Chemistry, Elsevier ScienceDirect, and Springer Nature. Each publication contains the full manuscript text, including experimental sections, results, and supplementary materials where available. This corpus represents a comprehensive cross-section of the FET sensor literature published over the past two decades, capturing both materials innovation and device engineering advances.

Importantly, the corpus is filtered to include only experimental studies (excluding computational simulations such as DFT or AIMD), and focuses specifically on concentration-detection applications spanning gas-phase, chemical, and biological FET sensors. Publications concerning non-concentration measurements (e.g., pressure sensing) are excluded from this dataset.

In that prior study [Ferreira et al. \(2025\)](#), domain experts manually curated structured metadata from these publications, establishing a high-quality ground-truth dataset. Each paper was annotated to extract eight core performance and operational parameters that characterize FET sensor functionality. These expert annotations serve as the gold standard for evaluating automated extraction methods in the present work.

Structured Information Extraction Task The primary objective of this work is to develop an automated pipeline capable of extracting structured sensor metadata from unstructured scientific text with fidelity approaching human expert performance. Specifically, given the full text of a scientific publication, the system must generate a structured JSON output containing the following fields:

- `sensor_type`: Classification of the sensor modality (gas, bio, or liquid)
- `detect_target`: Chemical or biological species being detected (e.g., ammonia, glucose, DNA)
- `lower_detection_limit`: Minimum detectable concentration with units (e.g., 5 ppb, 1 nM)
- `upper_detection_limit`: Maximum measurable concentration before saturation
- `probe_material`: Active sensing material or functionalization layer (e.g., metal oxide, antibody, polymer)
- `test_operating_temperature`: Experimental temperature in degrees Celsius
- `pH_value`: Solution pH for liquid/bio sensors, or -1 for gas-phase measurements
- `test_medium`: Environmental context of sensing experiments (e.g., air, phosphate buffer, serum)

This extraction task is challenging due to several factors: (1) high variability in reporting conventions across journals and research groups, (2) implicit information requiring domain knowledge inference (e.g., room temperature defaults, pH assumptions), (3) multi-scale unit conversions (ppb, ppm, molarity), and (4) disambiguation when multiple sensor variants are presented in a single publication. Moreover, detection limits are often reported indirectly through figures or cited without explicit numerical statements, requiring contextual reasoning beyond simple pattern matching.

Data Partitioning To ensure robust evaluation and generalization assessment, we employ a three-fold cross-validation strategy. The 844 publications are partitioned into three non-overlapping folds, each maintaining representative coverage of sensor types and target analytes. For each fold, the data is divided into:

- Training set: 583 publications (approximately 69 percent) for prompt optimization and model development

- 2160 • Test set: 261 publications (approximately 31 percent) held out for final performance eval-
2161 uation
2162

2163 The fold assignments are deterministic and fixed across all experiments to enable fair comparison
2164 between different optimization strategies. Critically, the test set remains completely isolated from
2165 the training process; no information from test set papers influences prompt refinement or model
2166 selection. The development set is available but not actively utilized in the current study, as our
2167 TextGrad-based optimization operates exclusively on training set feedback.

2168 All data splits preserve the original expert annotations, ensuring that each publication in the corpus
2169 has a corresponding ground-truth JSON record. This pairing enables automated metric computa-
2170 tion by comparing system-generated outputs against expert-curated references across multiple di-
2171 mensions: lexical overlap (BLEU, ROUGE), semantic similarity (BERTScore), and structured field
2172 accuracy (Exact Match, Jaccard Index).

2173
2174 **Extended Field Extraction** Beyond the eight core fields annotated in the prior work [Ferreira](#)
2175 [et al. \(2025\)](#), we extend the extraction scope to four additional field categories encompassing 20
2176 supplementary parameters:

2177
2178 **Electrode Architecture (4 fields)** Device electrode configuration including gate, source, and drain
2179 electrode materials, along with structure design type (e.g., planar, vertical, coplanar, back-gated, top-
2180 gated, interdigitated).

2181
2182 **Material Layer Composition (7 fields)** Layer-by-layer device structure comprising substrate ma-
2183 terial and thickness, channel (active layer) material, dielectric layer material and thickness, surface
2184 functionalization molecules, and structure dimensionality (0D/1D/2D/3D nanostructure classifica-
2185 tion).

2186
2187 **Sensitivity and Response Characteristics (4 fields)** Dynamic performance metrics including re-
2188 sponse time, recovery time, and sensitivity definition (numerator representing signal change such
2189 as ΔR , ΔI , or ΔG , and denominator representing reference value such as baseline resistance or
2190 analyte concentration).

2191
2192 **Synthesis and Thermal Processing (5 fields)** Fabrication conditions encompassing annealing
2193 temperature, duration, and atmosphere (air, N₂, Ar, O₂, vacuum, H₂), as well as hydrothermal
2194 synthesis temperature and duration.

2195 For these extended fields, expert annotations are performed on a smaller corpus of 171 publications
2196 (approximately 20 percent of the core dataset size), reflecting a deliberate strategy to minimize
2197 human labeling effort while maintaining sufficient data for prompt optimization.

2198 Our optimization workflow proceeds in two stages. First, the eight optimization modes (M1-M8) are
2199 systematically evaluated on the core dataset (844 papers) to identify the top-performing strategies.
2200 The three best-performing strategies are then applied to optimize prompts for all five field categories
2201 (core plus four extended), yielding 15 optimized prompts in total. In the final deployment phase,
2202 these prompts are applied to the full corpus of 1,686 publications. For each field, human experts
2203 review the outputs from the three strategy variants and perform lightweight curation to produce the
2204 final structured database. This ensemble approach, combining multiple optimization strategies with
2205 expert oversight, ensures robust extraction quality while leveraging the complementary strengths of
2206 different optimization modes.

2207
2208 **Problem Formulation** The central research question addressed in this phase is: Can an au-
2209 tonomous prompt optimization framework iteratively refine a natural language prompt to achieve
2210 expert-level extraction performance on complex scientific information extraction tasks, without
2211 manual trial-and-error? We begin with a minimal human-authored prompt (18 lines, providing only
2212 a JSON schema template) and seek to evolve it through data-driven optimization into a compre-
2213 hensive extraction guideline that encodes domain conventions, edge case handling, and implicit
reasoning strategies. Success is measured by comparing the optimized prompt against the baseline
across six complementary metrics computed on the held-out test set, with particular emphasis on

the Committee Score (average of BERTScore, ROUGE-1, and METEOR) as the primary quality indicator.

C.1.2 AUTOMATIC PROMPT OPTIMIZATION VIA TEXTGRAD

Overview Traditional prompt engineering for large language models relies on manual trial-and-error refinement, which is labor-intensive, subjective, and difficult to scale across diverse tasks. To address this limitation and enable autonomous prompt refinement, we adopt TextGrad (Yuksekgonul et al. (2025b)), a gradient-based optimization framework that treats natural language prompts as differentiable parameters. Analogous to how neural network training computes numerical gradients to update model weights, TextGrad computes *textual gradients*—natural language critiques generated by an LLM that describe how the prompt should be modified to improve task performance. This formulation conceptualizes the LLM as a differentiable computational graph where the prompt serves as a learnable input, the extraction task defines the loss function, and the critique model provides the optimization signal, enabling principled iterative refinement without manual intervention.

In this work, we instantiate TextGrad through a two-model architecture specifically designed for scientific information extraction tasks. Our implementation extends the original TextGrad framework by introducing minibatch training, hierarchical evaluation strategies, and multi-round iterative refinement to handle the complexity and variability inherent in real-world scientific literature.

Two-Model Architecture Our optimization pipeline employs a division of labor between two distinct language model components, each serving a specialized function in the prompt refinement process.

Inference Model The Inference Model serves as the production system responsible for performing the actual information extraction task. Given a scientific publication’s full text and a candidate prompt, the Inference Model generates structured JSON output containing the eight target fields described in Section 2.1. We evaluate multiple open-source models in this role, primarily focusing on DeepSeek-R1 variants (14B and 70B parameter versions) due to their strong reasoning capabilities and local deployment feasibility.

The Inference Model operates under inference mode only and does not undergo any fine-tuning or parameter updates. Its behavior is entirely controlled by the prompt text, making prompt quality the sole determinant of extraction performance. This constraint aligns with practical deployment scenarios where model retraining is infeasible due to computational costs or proprietary model restrictions.

Critique Model The Critique Model functions as a meta-optimizer that analyzes the performance gap between Inference Model outputs and expert annotations. Operating at a higher level of abstraction, the Critique Model receives as input: (1) the current prompt text, (2) a minibatch of scientific papers, (3) the corresponding Inference Model outputs for those papers, and (4) the ground-truth expert annotations. Its task is to identify systematic deficiencies in the current prompt and generate textual feedback (the “gradient”) that guides prompt revision.

We explore several model choices for this critique role, including DeepSeek-R1 (14B/70B), Qwen3 (235B parameters), and proprietary models such as GPT-o4-mini-high. The choice of Critique Model represents a key design decision, as stronger reasoning capabilities in this component can lead to more effective prompt optimization, albeit at higher computational cost.

Critically, the Critique Model does not directly rewrite the prompt. Instead, it produces an intermediate representation: a natural language instruction enumerating specific issues to address (e.g., “The prompt fails to handle pH conversion from logarithmic scale” or “Detection limit extraction ignores P/Pv ratio formats”). This separation ensures that prompt modifications remain interpretable and traceable.

Updater Component Following critique generation, a third component (the Updater, also implemented as an LLM call) consumes the critique’s textual gradient alongside the current prompt and produces a revised prompt (Eq. 2). The textual gradient identifies a small number of targeted deficiencies relative to expert annotations (typically 1–3 per round). To encourage incremental re-

2268 refinement rather than drastic rewrites, the Updater instruction limits the scope of each revision to a
 2269 modest fraction of the prompt tokens (e.g., up to roughly 25%), preventing wholesale replacement
 2270 that might discard functional aspects of the existing prompt. This is a soft constraint conveyed
 2271 through the Updater’s instruction rather than a guaranteed token-level projection, analogous in spirit
 2272 to gradient descent’s incremental update principle adapted to the discrete space of natural language.
 2273

2274 **Single-Round Optimization Workflow** Algorithm 1 formalizes the optimization procedure for
 2275 a single training round. The process begins by sampling a minibatch of papers from the training
 2276 set. For each paper, the Inference Model generates an extraction attempt using the current prompt
 2277 (Phase 1), producing a fixed set of baseline outputs that will be reused throughout the optimization
 2278 iterations.

2279 The Critique Model then analyzes this minibatch-level performance data, identifying systematic
 2280 prompt deficiencies that manifest across multiple papers. This aggregation over a minibatch (rather
 2281 than optimizing on individual papers) is crucial for learning generalizable prompt improvements
 2282 rather than overfitting to idiosyncrasies of single documents.
 2283

2284 The textual gradient produced by the Critique Model serves as input to the Updater, which performs
 2285 a constrained revision of the prompt (Phase 2). Critically, both optimization iterations analyze the
 2286 same fixed baseline outputs O_0 rather than regenerating outputs after each prompt revision. Op-
 2287 tionally, quantitative metrics computed from O_0 can be included in the critique prompt to guide the
 2288 analysis, though by default the Critique Model operates on qualitative comparison alone. This de-
 2289 sign reduces computational cost by avoiding redundant Inference Model calls during the refinement
 2290 loop, at the trade-off of operating on potentially outdated performance signals.

2291 After completing both iterations, the optimized prompts P_1 and P_2 are applied to the minibatch
 2292 to generate fresh outputs O_1 and O_2 (Phase 3). Finally, a winner selection mechanism determines
 2293 which prompt variant will propagate to the next training round based on these newly generated
 2294 outputs (Phase 4). This selection can be based purely on automated metric scores (hard evaluation)
 2295 or incorporate LLM-based pairwise comparison (soft evaluation), as discussed in Section 2.3. When
 2296 enforced child replacement is enabled, an optimized prompt is selected even if the initial prompt
 2297 achieved higher scores.

2298 **Algorithm 1** TextGrad Single-Round Optimization

2299 **Require:** Training set $\mathcal{D}_{\text{train}}$, current prompt P_0 , minibatch size B , iterations K
 2300 **Ensure:** Optimized prompt P^* for next round
 2301 1: Sample minibatch $\mathcal{M} \leftarrow \text{RandomSample}(\mathcal{D}_{\text{train}}, B)$
 2302 2: ▷ Phase 1: Generate baseline outputs with initial prompt
 2303 3: $\mathcal{O}_0 \leftarrow \{\text{InferenceModel}(P_0, \text{paper}) \mid \text{paper} \in \mathcal{M}\}$
 2304 4:
 2305 5: ▷ Phase 2: Iterative prompt refinement via textual gradients
 2306 6: Initialize $P \leftarrow P_0$
 2307 7: **for** $k = 1$ to K **do**
 2308 8: $I_k \leftarrow \text{CritiqueModel}(P, \mathcal{M}, \mathcal{O}_0, \mathcal{M}_{\text{expert}})$ ▷ Analyze \mathcal{O}_0 vs expert
 2309 9: $P_k \leftarrow \text{Updater}(P, I_k)$ ▷ Apply textual gradient
 2310 10: Save P_k to disk
 2311 11: $P \leftarrow P_k$ ▷ Chain for next iteration
 2312 12: **end for**
 2313 13:
 2314 14: ▷ Phase 3: Re-evaluate with optimized prompts
 2315 15: **for** $k = 1$ to K **do**
 2316 16: $\mathcal{O}_k \leftarrow \{\text{InferenceModel}(P_k, \text{paper}) \mid \text{paper} \in \mathcal{M}\}$
 2317 17: **end for**
 2318 18:
 2319 19: ▷ Phase 4: Compute metrics and select winner
 2320 20: $\mathcal{S} \leftarrow \text{EvaluateMetrics}(\{\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_K\}, \mathcal{M}_{\text{expert}})$
 2321 21: $P^* \leftarrow \text{SelectWinner}(\{P_0, P_1, \dots, P_K\}, \mathcal{S}, \text{duel_mode})$
 2322 22: **return** P^*

2322 **Multi-Round Training Strategy** The single-round optimization described above is embedded
 2323 within an outer loop that spans multiple training rounds (typically 20-100 rounds depending on
 2324 the optimization mode). Each round operates on a randomly sampled minibatch from the training
 2325 set, ensuring that the prompt is exposed to diverse examples over time rather than overfitting to a
 2326 fixed subset. Note that the sampling is performed without tracking or deduplication across rounds,
 2327 meaning that individual papers may appear in multiple minibatches over the course of training. This
 2328 design prioritizes simplicity and reproducibility (via fixed random seeds) over strict non-replacement
 2329 sampling.

2330 The winner prompt from round n serves as the initialization P_0 for round $n + 1$, creating a trajectory
 2331 of progressively refined prompts. This multi-round strategy allows the optimization to accumu-
 2332 late improvements incrementally, with early rounds addressing coarse-grained issues (e.g., output
 2333 format compliance) and later rounds refining subtle edge cases (e.g., unit conversion conventions,
 2334 disambiguation strategies).

2335 Importantly, the training set is never exhausted; with 583 training papers and a minibatch size of 3,
 2336 a 20-round optimization exposes the system to only 60 papers (approximately 10 percent of avail-
 2337 able data). This deliberate undersampling reduces computational cost while maintaining sufficient
 2338 diversity for generalization. The held-out test set (218 papers) provides an unbiased measure of the
 2339 final prompt’s performance on unseen data.

2340

2341 C.1.3 OPTIMIZATION MODES AND WINNER SELECTION STRATEGIES

2342

2343 A key contribution of this work is the systematic exploration of different optimization strategies
 2344 through a modular framework. Rather than committing to a single winner selection mechanism,
 2345 we define eight distinct optimization modes (M1-M8) that combine complementary evaluation ap-
 2346 proaches. This design enables empirical comparison of pure metric-driven optimization versus
 2347 LLM-assisted qualitative judgment, and allows us to assess the value of hierarchical competition
 2348 structures.

2349 These modes represent different combinations of configurable components: whether to include hard
 2350 score feedback in critique prompts, whether to employ LLM-based judging for winner selection,
 2351 and whether to enforce child replacement regardless of performance. Table S2 summarizes these
 2352 configurations as a reference taxonomy for the experimental evaluation.

2353

2354 **Framework Components** Our optimization modes are constructed from four orthogonal design
 2355 dimensions:

2356

2357 **Winner Selection Basis** The mechanism for determining which prompt variant advances to the
 2358 next training round. Two options are available:

- 2359 • **Hard Score:** Winner determined solely by automated metric performance via the Commit-
 2360 tee Score (Eq. 3), which averages over available metrics (default: BERTScore, ROUGE-1,
 2361 METEOR; alternatively “all”: BLEU, ROUGE-1, METEOR, BERTScore, ExactMatch,
 2362 Jaccard). This approach is deterministic, computationally efficient, and directly optimizes
 2363 the target evaluation function.
- 2364 • **Duel Result:** Winner determined by LLM-based pairwise comparison, where a judge model
 2365 evaluates which prompt produces outputs closer to expert quality. This approach captures
 2366 nuanced quality dimensions not fully reflected in lexical metrics.

2367

2368 **Duel Mode** The competition structure for comparing prompt variants within a training round:

2369

- 2370 • **No Duel:** Only automated metrics are computed; no LLM judging is performed. Fastest
 2371 option, suitable for large-scale experiments.
- 2372 • **Child-Only Duel:** The two optimized prompts (iteration-1 vs iteration-2) compete in pair-
 2373 wise LLM judging to determine a child champion. When combined with hard score winner
 2374 basis, this child champion then competes against the initial prompt using automated met-
 2375 rics. When combined with duel result winner basis, the LLM-selected child champion
 becomes the final winner directly.

- 2376
- 2377
- 2378
- 2379
- 2380
- 2381
- Hierarchical Duel: A two-stage tournament where (1) iteration-1 competes against iteration-2 via LLM judging to determine a child champion, then (2) the child champion competes against the initial prompt (parent) in a second LLM-judged comparison. This structure tests whether optimization genuinely improves over the baseline through qualitative assessment at both stages.

2382 **Feedback Inclusion** Whether to provide quantitative performance metrics to the Critique Model during gradient generation:

2383

- 2384
- 2385
- 2386
- 2387
- 2388
- 2389
- 2390
- 2391
- 2392
- Excluded (default): The Critique Model receives only the textual outputs and expert annotations, forcing it to identify deficiencies through qualitative comparison alone. This approach mirrors the original TextGrad formulation and avoids potential metric hacking.
 - Included: The Critique Model additionally receives aggregate hard scores for the current prompt (e.g., average BERTScore, ROUGE-1, METEOR, and Committee Score across the minibatch). This summary-level feedback can help prioritize improvement areas but may bias the optimization toward metric-driven refinements rather than semantic quality improvements.

2393 **Child Replacement Policy** Whether to enforce selection of an optimized prompt regardless of performance:

2394

2395

- 2396
- 2397
- 2398
- 2399
- 2400
- 2401
- Flexible: The initial prompt can win if it outperforms both optimized variants. This conservative approach prevents performance regression.
 - Enforced: One of the optimized prompts (iteration-1 or iteration-2) must be selected even if the initial prompt scored higher. This aggressive strategy forces exploration and prevents optimization stagnation.

2402 **Mode Definitions** Table S2 enumerates the eight optimization modes investigated in this study. The modes are ordered by increasing complexity and computational cost.

2403

2404

2405 Table S2: Optimization mode configurations. Hard score refers to automated metrics (BERTScore, ROUGE, etc.), while duel result indicates LLM-based pairwise judgment.

2406

2407

2408

Mode	Winner Basis	Duel Mode	Feedback	Child Enforced
M1	Hard Score	No Duel	Excluded	No
M2	Hard Score	No Duel	Included	No
M3	Hard Score	Child-Only	Excluded	No
M4	Duel Result	Child-Only	Excluded	No
M5	Hard Score	Child-Only	Included	No
M6	Duel Result	Hierarchical	Excluded	No
M7	Duel Result	Hierarchical	Excluded	Yes
M8	Duel Result	Hierarchical	Included	Yes

2416

2417 Mode Rationale and Expected Behaviors

2418

2419 **M1: Hard Score Baseline** The simplest configuration, serving as a baseline for pure metric-driven optimization. The Critique Model operates in a "black box" setting with no quantitative feedback, and winner selection relies entirely on automated metrics. This mode is fastest and most reproducible, but may struggle with quality dimensions poorly captured by lexical similarity measures.

2420

2421

2422

2423

2424 **M2: Metric-Aware Critique** Extends M1 by providing hard scores to the Critique Model. This allows the critique to diagnose specific metric weaknesses (e.g., "ROUGE is low due to missing synonyms") but risks overfitting to metric idiosyncrasies rather than true semantic quality.

2425

2426

2427

2428 **M3-M5: Child Competition with Hard Score Selection** These modes introduce LLM-based judging for comparing the two optimized prompts (iteration-1 vs iteration-2), with final winner selection based on hard scores. In M3, the LLM first determines which child prompt is superior,

2429

then this child champion competes against the parent prompt using automated metrics. M5 extends M3 by providing aggregate performance metrics to the Critique Model during gradient generation, enabling metric-aware prompt refinement. Both modes leverage LLM judgment for relative child comparison while relying on objective metrics for the critical parent-vs-child decision.

M4: Child Competition with LLM Selection M4 uses LLM-based judging throughout the selection process. The two optimized prompts compete via pairwise LLM comparison, and the winning child is directly selected as the round winner without further comparison against the parent prompt. This aggressive strategy fully trusts LLM judgment to identify improvements, bypassing metric-based validation against the baseline.

M6: Hierarchical Tournament Implements a two-stage competition where (1) the optimized prompts first compete via LLM judging to select a child champion, then (2) the child champion faces the parent prompt in a second LLM-judged comparison. This structure mirrors evolutionary selection pressure, ensuring that optimized variants must demonstrably outperform the baseline through qualitative assessment to propagate.

M7: Forced Exploration Enforces child selection even when the parent outperforms, preventing premature convergence to local optima. This mode is appropriate when early-round regressions are tolerable in pursuit of eventual breakthroughs, or when hard metrics are known to be unreliable guides.

M8: Full System Combines all enhancement components: hierarchical dueling, LLM-based winner selection, metric-aware critique, and enforced exploration. This mode represents maximum computational investment and is expected to achieve the most aggressive optimization, though at risk of instability or overfitting to the critique model’s biases.

Computational Complexity The eight modes exhibit substantial variation in computational cost. Modes M1-M2 require only Inference Model inference and metric calculation, while M8 additionally invokes the Critique Model, Updater, and multiple LLM judge calls per training round. For a 20-round optimization on fold 1 (60 papers total), M1 requires approximately 120 LLM calls (60 initial + 60 re-evaluation after two iterations), whereas M8 may require over 500 calls when accounting for critique, update, and judging operations.

This cost-performance trade-off is a central empirical question: Does the additional reasoning provided by LLM-based judging and metric-aware critique justify the 4-5x increase in computational expense? Our experimental results (Section 3) provide evidence to inform this design choice for practical deployments.

C.1.4 EVALUATION METRICS

Assessing the quality of automatically extracted information requires comparing system-generated outputs against expert-curated ground truth across multiple dimensions. No single metric can fully capture extraction fidelity, as structured scientific information encompasses lexical accuracy, semantic faithfulness, and field-level completeness. We therefore employ a complementary suite of six automated metrics that collectively evaluate different aspects of extraction performance.

Semantic Similarity: BERTScore [Zhang et al. \(2020\)](#) BERTScore measures semantic similarity by computing contextual embeddings for tokens in both the generated and reference texts, then finding optimal alignment between token pairs using cosine similarity. Unlike surface-form metrics, BERTScore captures paraphrases and synonymous expressions through its pre-trained language model backbone (DistilBERT-base-uncased in our implementation).

For each token in the candidate output, the metric identifies the most similar token in the reference based on embedding distance, and vice versa. Precision, recall, and F1 scores are computed from these alignments, with F1 serving as the primary BERTScore value. This metric is particularly valuable for scientific text where authors may describe the same concept using varied terminology (e.g., "detection limit" versus "sensitivity threshold").

2484 The key advantage of BERTScore is its robustness to surface-form variation while maintaining inter-
2485 interpretability through token-level alignments. However, it may occasionally conflate semantically
2486 distinct technical terms with similar contextual usage patterns, and its reliance on a fixed embedding
2487 model means it cannot adapt to domain-specific terminology beyond its training corpus.

2488
2489 **Lexical Overlap: ROUGE-1** [Lin \(2004\)](#) ROUGE-1 (Recall-Oriented Understudy for Gisting
2490 Evaluation) measures unigram overlap between generated and reference texts. Originally designed
2491 for summarization evaluation, ROUGE-1 quantifies how much of the reference content is captured
2492 in the system output, emphasizing recall over precision.

2493 The metric operates by stemming all words to their root forms (using the Porter stemmer), then
2494 computing the ratio of overlapping unigrams to total unigrams in the reference. ROUGE-1 is par-
2495 ticularly sensitive to completeness: a system that omits critical fields or values will incur steep
2496 penalties. However, it is agnostic to word order and semantic nuance, treating "5 ppb" and "5 parts
2497 per billion" as entirely distinct despite equivalent meaning.

2498 We report the ROUGE-1 F1 score, which balances recall against precision, preventing trivial opti-
2499 mization strategies that maximize overlap by copying entire passages verbatim.

2500
2501 **Alignment-Based Evaluation: METEOR** [Banerjee & Lavie \(2005\)](#) METEOR (Metric for Eval-
2502 uation of Translation with Explicit ORdering) extends simple n-gram matching by incorporating
2503 stemming, synonymy, and paraphrase recognition. The metric aligns words between candidate and
2504 reference texts using multiple matching stages: exact match, stem match, and synonym match (based
2505 on WordNet).

2506 After establishing alignments, METEOR computes precision and recall, then combines them into a
2507 harmonic mean weighted toward recall. Additionally, METEOR applies a fragmentation penalty that
2508 rewards contiguous matches over scattered alignments, implicitly capturing some notion of fluency
2509 and coherence.

2510 For structured extraction tasks, METEOR's synonym awareness is particularly valuable when deal-
2511 ing with chemical nomenclature or equivalent technical expressions. For instance, "phosphate
2512 buffered saline" and "PBS" would receive partial credit despite sharing no surface tokens. This
2513 makes METEOR more lenient than BLEU while remaining more conservative than pure semantic
2514 embeddings.

2515
2516 **N-Gram Precision: BLEU** [Papineni et al. \(2002\)](#) BLEU (Bilingual Evaluation Understudy)
2517 measures precision of n-gram matches between generated and reference texts, with a brevity penalty
2518 to discourage pathologically short outputs. Originally designed for machine translation, BLEU em-
2519 phasizes exactness over recall: a system that produces highly accurate but incomplete extractions
2520 can still achieve respectable BLEU scores.

2521 We employ the smoothing variant (method 1) to handle cases where higher-order n-grams may have
2522 zero matches, preventing undefined scores on short outputs. BLEU's primary limitation for extrac-
2523 tion tasks is its insensitivity to semantic equivalence; functionally correct variations in phrasing or
2524 unit representation receive no credit unless they exhibit surface-level overlap.

2525 Despite these limitations, BLEU remains a useful signal for detecting prompt refinements that im-
2526 prove format consistency and terminology standardization.

2527
2528 **Structured Field Accuracy: Exact Match** Exact Match evaluates field-level correspondence in
2529 the structured JSON outputs. For each of the eight target fields (sensor type, detect target, detec-
2530 tion limits, etc.), we compare the extracted value against the expert annotation and compute the
2531 percentage of fields with character-exact agreement.

2532 This metric is unforgiving: "5.0 ppb" and "5 ppb" are considered distinct, as are "ammonia" and
2533 "NH3". Exact Match therefore captures the most stringent notion of extraction correctness, reward-
2534 ing systems that perfectly replicate expert formatting conventions and terminology choices.

2535
2536 The metric's harshness makes it a strong indicator of prompt refinement quality. Improvements in
2537 Exact Match signal that the system has learned to standardize units, resolve abbreviations consis-
tently, and adopt expert-preferred field values rather than extracting raw text snippets.

Set Similarity: Jaccard Index The Jaccard Index measures token-level set similarity as the ratio of intersecting tokens to the union of tokens across generated and reference outputs. After tokenization (splitting on whitespace and punctuation), we compute:

$$\text{Jaccard} = \frac{|\text{tokens}_{\text{gen}} \cap \text{tokens}_{\text{ref}}|}{|\text{tokens}_{\text{gen}} \cup \text{tokens}_{\text{ref}}|}$$

This metric is order-agnostic and emphasizes coverage: outputs that mention the same entities and values as the reference, even in different arrangements, receive high scores. Jaccard is less sensitive to verbosity than ROUGE (since it uses union normalization) and less sensitive to phrasing than BLEU (since it ignores n-gram structure).

For multi-field structured outputs, Jaccard provides a complementary signal to Exact Match. While Exact Match requires perfect field-level agreement, Jaccard rewards partial matches and gives credit for extracting most of the relevant information even when formatting details differ.

Metric Complementarity and Composite Scoring The six base metrics above capture orthogonal quality dimensions: BERTScore for semantic fidelity, ROUGE for completeness, METEOR for flexible lexical alignment, BLEU for precision, Exact Match for structural correctness, and Jaccard for entity coverage. No single metric is universally superior; each exhibits distinct failure modes and biases.

To support winner selection decisions, we additionally compute a Committee Score by averaging BERTScore, ROUGE-1, and METEOR. This composite metric (yielding seven total scores per prompt variant) provides a balanced summary emphasizing the most reliable evaluation dimensions. During optimization, these quantitative scores may optionally be provided to the Critique Model to inform its analysis, though by default the critique operates on qualitative output comparison alone.

All final experimental results report the complete metric profile (six base metrics plus committee score), enabling comprehensive assessment of prompt quality improvements across multiple evaluation axes. This multi-metric approach guards against overfitting to any single measure and provides a more robust characterization of extraction fidelity.

C.1.5 EXPERIMENTAL PROTOCOL

Training Phase The prompt optimization process operates over multiple training rounds; we sweep {20, 40, 60, 100} rounds for every configuration rather than tying round count to a specific mode. Each round executes the single-round workflow described in Algorithm 1: a minibatch of 3 papers is randomly sampled from the training set, the Inference Model generates outputs using the current prompt, the Critique Model analyzes performance gaps, and the Updater produces a revised prompt. Winner selection (via hard scores, LLM judging, or hierarchical dueling) determines which prompt variant propagates to the subsequent round.

This multi-round strategy ensures exposure to diverse examples over time while avoiding exhaustive iteration over the entire training set. With 583 training papers available and a minibatch size of 3, the round-count sweep spans 60–300 unique paper encounters per trajectory (with possible repeats due to random sampling). This deliberate undersampling reduces computational cost while maintaining sufficient diversity for generalization.

Random sampling is controlled by a fixed seed per experimental run to ensure reproducibility. Each training round operates on a freshly sampled minibatch, preventing overfitting to a static subset of examples and allowing the optimization to encounter different edge cases and reporting conventions throughout the training trajectory.

Test Phase Evaluation Upon completion of multi-round training, the final optimized prompt is evaluated on the held-out test set (261 papers per fold). The Inference Model processes each test paper using both the original human-written prompt and the optimized prompt, generating paired outputs for direct comparison.

For each test paper, we compute the six-metric evaluation suite (BERTScore, ROUGE-1, METEOR, BLEU, Exact Match, Jaccard) by comparing system outputs against expert annotations. Aggregate

2592 statistics (mean, standard deviation, minimum, maximum) are reported across all test papers for
2593 each metric, along with per-metric improvement ratios between optimized and baseline prompts.

2594 Critically, no information from test set papers influences the optimization process. Test set DOIs,
2595 full texts, and expert annotations remain isolated from training, ensuring that reported performance
2596 reflects genuine generalization to unseen scientific literature rather than memorization of training
2597 examples.

2598
2599 **Experimental Design and Replication** To assess optimization robustness and account for
2600 stochastic variation, we employ a rigorous replication strategy combining multiple independent runs
2601 with cross-validation:

2602
2603 **Independent Replications** For each configuration (optimization mode, model choice, hyperpa-
2604 rameters), we conduct 5 independent rollouts starting from the same human-authored initial prompt.
2605 Seeds are fixed for reproducibility, but minibatch sampling order differs across rollouts, yielding dis-
2606 tinct optimization trajectories even under identical settings. Final test set performance is aggregated
2607 across these 5 runs, with mean and standard deviation reported to characterize typical performance
2608 and run-to-run variability.

2609
2610 **Cross-Validation** The 3-fold data partitioning enables assessment of methodology-level variation
2611 beyond single-fold idiosyncrasies. By repeating experiments across all three folds, we obtain per-
2612 formance estimates that account for dataset composition effects. Results tables report both per-fold
2613 performance (to demonstrate consistency) and cross-fold aggregates (to summarize overall effec-
2614 tiveness).

2615 This factorial design (5 replications \times 3 folds = 15 total runs per configuration) provides statistical
2616 rigor for comparing optimization modes and model choices.

2617 C.2 METHOD AND TECHNICAL DETAILS FOR “TWIN”

2618 C.2.1 MATERIAL DESCRIPTOR REPRESENTATION

2619
2620 Field-effect transistor (FET) sensors incorporate diverse material classes across functional compo-
2621 nents, including channel semiconductors, dielectric layers, electrodes, substrates, and surface func-
2622 tionalizations. To enable machine learning on these heterogeneous materials, we develop a unified
2623 descriptor framework that represents each material as a fixed-dimensional numerical vector com-
2624 prising 25 macroscopic properties and a 320-dimensional fingerprint embedding.

2625
2626 Building upon prior work on spiking graph neural networks for FET sensor modeling [Ferreira et al.](#)
2627 [\(2025\)](#), which demonstrates the effectiveness of material descriptors for cheminformatics and mate-
2628 rials informatics, we extend the framework to accommodate the broader scope of the present dataset.
2629 While the original implementation covered inorganic compounds, organic molecules, and polymers,
2630 the expanded corpus from the Text phase now includes biosensors with protein-based recognition
2631 elements and nucleic acid aptamers. Accordingly, we introduce two additional material categories—
2632 biomolecules (proteins) and nucleic acids (DNA/RNA)—leveraging state-of-the-art biological lan-
2633 guage models for their representation.

2634 This *Cross-Domain Material Fingerprinting* approach transforms free-text material names extracted
2635 during the Text phase into dense numerical vectors, substantially enriching the structured informa-
2636 tion available for the Twin phase. Each publication may yield one or more experimental records,
2637 and our encoding scheme ensures that the full chemical and biological diversity across all material-
2638 bearing fields is captured in a unified, model-ready format.

2639
2640 **Material Classification and Descriptor Sources** Materials extracted from FET sensor literature
2641 were classified into five categories based on chemical structure: inorganic compounds, organic
2642 molecules, polymers, biomolecules (proteins), and nucleic acids (DNA/RNA). Each category em-
2643 ploys domain-specific descriptors optimized for capturing physicochemical properties relevant to
2644 sensing performance (Table [S3](#)).

2645 For inorganic materials (Table [S4](#)), we retrieved formation energies, band gaps, elastic moduli, di-
electric constants, and effective masses from the Materials Project database via the OPTIMADE

Table S3: Material descriptor sources and feature dimensions by material category.

Category	Macroproperties (25D)	Fingerprint (320D)	Database
Inorganic	14 physical properties + 11 one-hot encoded (crystal system, electronic class)	MAGPIE composition embedding	Materials Project Jain et al. (2013)
Molecules	Electronic, polarity, molecular size, flexibility, and complexity descriptors	Morgan circular fingerprint	PubChem National Center for Biotechnology Information (2024) , ChEMBL European Bioinformatics Institute (2024)
Polymers	7 bulk properties (thermal, mechanical, electrical) + 18 monomer descriptors	Morgan fingerprint of repeat unit	PubChem National Center for Biotechnology Information (2024)
Biomolecules	Sequence length, physicochemical properties, amino acid composition, secondary structure	ESM-2 protein embedding	UniProt UniProt Consortium (2015) , ESM-2 Lin et al. (2023)
DNA/RNA	Sequence composition, GC content, melting temperature, thermodynamic stability	DNABERT-S embedding	DNABERT-S Zhou et al. (2025)

API. Crystal system (7 classes) and electronic classification (4 classes: metal, semiconductor, insulator, other) were one-hot encoded to complete the 25-dimensional macroproperty vector. The 320-dimensional fingerprint was derived from MAGPIE (Materials-Agnostic Platform for Informatics and Exploration) composition-based features.

Organic molecules (Table [S5](#)) and polymer repeat units (Table [S6](#)) are characterized using molecular descriptors from PubChem and ChEMBL, including topological polar surface area (TPSA), partition coefficients (XLogP), hydrogen bond donor/acceptor counts, and 3D pharmacophore features. Morgan circular fingerprints (radius 2, 256 bits) are computed using RDKit and zero-padded to 320 dimensions.

Protein descriptors (Table [S7](#)) are computed from amino acid sequences retrieved from UniProt, including isoelectric point, instability index, GRAVY hydrophobicity, and amino acid composition fractions. The 320-dimensional embedding is generated using the ESM-2 protein language model (650M parameters), with the mean-pooled representation extracted from the final hidden layer.

Nucleic acid sequences (Table [S8](#)) are characterized by length, GC content, purine/pyrimidine ratios, and predicted melting temperatures. Embeddings are computed using DNABERT-S, a species-aware DNA language model, yielding 256-dimensional vectors that are zero-padded to 320 dimensions. Table [S9](#) summarizes the fingerprint and embedding sources for all material categories.

Feature Aggregation Each FET sensor record contains up to 11 material-bearing fields: channel, dielectric layer, gate electrode, source electrode, drain electrode, substrate, probe material, surface functionalization, detection target, test medium, and annealing atmosphere. For each field, the corresponding material is mapped to its descriptor vector (25 + 320 = 345 dimensions). Composite materials (e.g., “zinc oxide/graphene”) are decomposed into individual components, and their descriptor vectors are averaged element-wise.

The complete feature representation for each sensor record is constructed by concatenating descriptors from all 11 material fields, yielding a $(11 \times 345) = 3795$ -dimensional material descriptor vector. Combined with 7 scalar device parameters (detection limits, response times, etc.) and categorical encodings, the final feature dimension is 3869.

C.2.2 HETEROGENEOUS GRAPH REPRESENTATION

To capture the complex multi-component architecture of FET sensors, we represent each device as a heterogeneous graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where nodes \mathcal{V} correspond to functional components and edges \mathcal{E} encode physical interactions between them. This representation directly mirrors the physical structure of FET sensors: discrete material components (channel, electrodes, dielectrics) interact

2700
2701
2702
2703
2704
2705
2706
2707
2708
2709
2710
2711
2712
2713
2714
2715
2716
2717
2718
2719
2720
2721
2722
2723
2724
2725
2726
2727
2728
2729
2730
2731
2732
2733
2734
2735
2736
2737
2738
2739
2740
2741
2742
2743
2744
2745
2746
2747
2748
2749
2750
2751
2752
2753

Table S4: Macroscopic properties for inorganic materials (25 dimensions).

#	Property	Unit	Description
<i>Thermodynamic & Electronic (14 dimensions)</i>			
1	formation_energy_per_atom	eV/atom	DFT formation energy from elemental states
2	energy_above_hull	eV/atom	Energy distance to convex hull (0 = stable)
3	band_gap	eV	Electronic band gap
4	density	g/cm ³	Mass density
5	epsilon_x	–	Dielectric constant (x-axis)
6	epsilon_y	–	Dielectric constant (y-axis)
7	epsilon_z	–	Dielectric constant (z-axis)
8	dielectric_total	–	Total static dielectric constant
9	k_vrh	GPa	Bulk modulus (Voigt-Reuss-Hill)
10	g_vrh	GPa	Shear modulus (Voigt-Reuss-Hill)
11	poisson_ratio	–	Poisson's ratio
12	me_avg	m ₀	Average electron effective mass
13	mh_avg	m ₀	Average hole effective mass
14	magnetization	μ _B /f.u.	Total magnetization
<i>Crystal System One-Hot (7 dimensions)</i>			
15	cs_cubic	0/1	Cubic crystal system
16	cs_hexagonal	0/1	Hexagonal crystal system
17	cs_orthorhombic	0/1	Orthorhombic crystal system
18	cs_tetragonal	0/1	Tetragonal crystal system
19	cs_trigonal	0/1	Trigonal crystal system
20	cs_monoclinic	0/1	Monoclinic crystal system
21	cs_triclinic	0/1	Triclinic crystal system
<i>Electronic Classification One-Hot (4 dimensions)</i>			
22	eg_metal	0/1	Metallic (band gap = 0)
23	eg_semiconductor	0/1	Semiconductor (0 < band gap < 3 eV)
24	eg_insulator	0/1	Insulator (band gap ≥ 3 eV)
25	eg_other	0/1	Other/unknown classification

2754
2755
2756
2757
2758
2759
2760
2761
2762
2763
2764
2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807

Table S5: Macroscopic properties for organic molecules (25 dimensions).

#	Property	Unit	Description
<i>Electronic Properties (5 dimensions)</i>			
1	Charge	e	Formal molecular charge
2	aromatic_rings	count	Number of aromatic ring systems
3	FeatureRingCount3D	count	3D aromatic ring pharmacophore features
4	FeatureCationCount3D	count	Cationic center count
5	FeatureAnionCount3D	count	Anionic center count
<i>Polarity & Surface Interactions (6 dimensions)</i>			
6	TPSA	Å ²	Topological polar surface area
7	XLogP	log units	Octanol-water partition coefficient
8	HBondDonorCount	count	Hydrogen bond donor count
9	HBondAcceptorCount	count	Hydrogen bond acceptor count
10	FeatureDonorCount3D	count	3D H-bond donor features
11	FeatureAcceptorCount3D	count	3D H-bond acceptor features
<i>Molecular Size & Shape (6 dimensions)</i>			
12	MolecularWeight	Da (g/mol)	Molecular weight
13	HeavyAtomCount	count	Non-hydrogen atom count
14	Volume3D	Å ³	Van der Waals volume
15	XStericQuadrupole3D	–	X-direction steric quadrupole moment
16	YStericQuadrupole3D	–	Y-direction steric quadrupole moment
17	ZStericQuadrupole3D	–	Z-direction steric quadrupole moment
<i>Flexibility & Dynamics (3 dimensions)</i>			
18	RotatableBondCount	count	Rotatable single bond count
19	EffectiveRotorCount3D	count	Effective rotor count in 3D
20	ConformerModelRMSD3D	Å	Conformer ensemble RMSD
<i>Complexity & Drug-likeness (5 dimensions)</i>			
21	Complexity	–	Bertz complexity score
22	FeatureHydrophobeCount3D	count	Hydrophobic pharmacophore features
23	FeatureCount3D	count	Total pharmacophore features
24	qed_weighted	0–1	Quantitative drug-likeness score
25	np_likeness_score	–5 to +5	Natural product likeness

Table S6: Macroscopic properties for polymers (25 dimensions).

#	Property	Unit	Description
<i>Bulk Polymer Properties (7 dimensions)</i>			
1	molar_mass_repeat.g_mol	g/mol	Repeat unit molecular weight
2	density.g_cm3	g/cm ³	Bulk density
3	Tg_C	°C	Glass transition temperature
4	Td_C	°C	Decomposition temperature
5	Tm_C	°C	Melting temperature
6	dielectric_constant	–	Relative permittivity
7	youngs_modulus.GPa	GPa	Elastic modulus
<i>Monomer Descriptors (18 dimensions)</i>			
8	Charge	e	Monomer formal charge
9	TPSA	Å ²	Polar surface area
10	HBondDonorCount	count	H-bond donors
11	HBondAcceptorCount	count	H-bond acceptors
12	MolecularWeight	Da	Monomer molecular weight
13	HeavyAtomCount	count	Non-H atom count
14	RotatableBondCount	count	Rotatable bonds
15	XLogP	log units	Lipophilicity
16	Complexity	–	Molecular complexity
17	FeatureRingCount3D	count	3D ring count
18	FeatureCationCount3D	count	Cationic features
19	FeatureAnionCount3D	count	Anionic features
20	FeatureDonorCount3D	count	3D H-bond donors
21	FeatureAcceptorCount3D	count	3D H-bond acceptors
22	Volume3D	Å ³	Monomer volume
23	EffectiveRotorCount3D	count	Effective rotors
24	FeatureHydrophobeCount3D	count	Hydrophobic features
25	aromatic_rings	count	Aromatic ring count

through well-defined mechanisms (carrier transport, capacitive coupling, surface chemistry) that determine sensing performance.

Node Types and Features We define 16 node types organized into three functional categories:

- **Device components** (9 nodes): channel, gate (top/bottom), dielectric layer (top/bottom), floating gate, source, drain, and substrate. These represent the transistor’s electrical architecture.
- **Sensing components** (5 nodes): surface functionalization, probe material, detection target, test medium, and electrolyte. These capture the biochemical sensing interface.
- **Process/condition** (2 nodes): annealing process parameters and operating conditions (temperature, pH, etc.).

Each node’s feature vector is constructed from the material descriptors described in Section C.2.1. For material-bearing nodes, we concatenate the 25-dimensional macroproperty vector with a 25-dimensional binary mask indicating feature availability, plus a scalar for the number of constituent materials (for composites), yielding a 51-dimensional node feature. The 320-dimensional fingerprints are processed separately through a dedicated neural branch (Section C.2.4).

For nodes with missing materials, we apply zero-masking: the feature vector is set to zeros while the mask indicates complete absence, allowing the model to distinguish between missing data and zero-valued properties.

Edge Types and Physical Semantics Rather than using a fully connected graph, we define six edge types that encode distinct physical interactions governing FET sensor operation. This physics-informed topology ensures that message passing in the graph neural network follows actual charge transport, electrostatic coupling, and chemical interaction pathways:

2862
2863
2864
2865
2866
2867
2868
2869
2870
2871
2872
2873
2874
2875
2876
2877
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915

Table S7: Macroscopic properties for biomolecules/proteins (25 dimensions).

#	Property	Unit	Description
<i>Sequence Properties (3 dimensions)</i>			
1	sequence_length	aa	Amino acid count
2	molecular_weight_kDa	kDa	Protein molecular weight
3	isoelectric_point	pH	Isoelectric point (pI)
<i>Physicochemical Properties (5 dimensions)</i>			
4	aromaticity	0–1	Aromatic residue fraction (Phe, Trp, Tyr)
5	instability_index	–	Protein stability index (>40 = unstable)
6	gravy	–	Grand average hydropathicity
7	charge_at_pH7	e	Net charge at pH 7.0
8	charge_at_pH5	e	Net charge at pH 5.0
<i>Amino Acid Composition (9 dimensions)</i>			
9	aromatic_fraction	%	Phe, Tyr, Trp fraction
10	aliphatic_fraction	%	Ala, Val, Leu, Ile fraction
11	polar_fraction	%	Ser, Thr, Asn, Gln fraction
12	charged_fraction	%	Asp, Glu, Lys, Arg fraction
13	positive_fraction	%	Lys, Arg, His fraction
14	negative_fraction	%	Asp, Glu fraction
15	cysteine_count	count	Disulfide-forming residues
16	proline_count	count	Structure-breaking residues
17	glycine_fraction	%	Flexible residue fraction
<i>Secondary Structure (3 dimensions)</i>			
18	helix_fraction	0–1	α -helix propensity
19	turn_fraction	0–1	β -turn propensity
20	sheet_fraction	0–1	β -sheet propensity
<i>Spectroscopic & Physical (5 dimensions)</i>			
21	extinction_coefficient_reduced	$M^{-1}cm^{-1}$	Extinction coeff. (reduced Cys)
22	extinction_coefficient_oxidized	$M^{-1}cm^{-1}$	Extinction coeff. (oxidized Cys)
23	average_flexibility	–	B-factor derived flexibility
24	tiny_fraction	%	Gly, Ala, Ser, Cys, Thr fraction
25	large_fraction	%	Phe, Ile, Lys, Leu, Met, Arg, Trp, Tyr fraction

2916
2917
2918
2919
2920
2921
2922
2923
2924
2925
2926
2927
2928
2929
2930
2931
2932
2933
2934
2935
2936
2937
2938
2939
2940
2941
2942
2943
2944
2945
2946
2947
2948
2949
2950
2951
2952
2953
2954
2955
2956
2957
2958
2959
2960
2961
2962
2963
2964
2965
2966
2967
2968
2969

Table S8: Macroscopic properties for DNA/RNA sequences (25 dimensions).

#	Property	Unit	Description
<i>Sequence Composition (9 dimensions)</i>			
1	length	nt	Nucleotide count
2	gc_content	%	Guanine + Cytosine fraction
3	at_content	%	Adenine + Thymine/Uracil fraction
4	a_count	count	Adenine count
5	c_count	count	Cytosine count
6	g_count	count	Guanine count
7	t_count	count	Thymine count
8	u_count	count	Uracil count (RNA only)
9	purine_pyrimidine_ratio	-	(A+G)/(C+T/U) ratio
<i>Thermodynamic Properties (3 dimensions)</i>			
10	tm_celsius	°C	Melting temperature
11	delta_g_kcal_mol	kcal/mol	Folding free energy
12	complexity	-	Sequence complexity index
<i>Sequence Features (4 dimensions)</i>			
13	longest_homopolymer	nt	Longest homopolymer run
14	cpg_count	count	CpG dinucleotide count
15	gc_skew	-	(G-C)/(G+C) strand asymmetry
16	at_skew	-	(A-T)/(A+T) strand asymmetry
<i>Dinucleotide Frequencies (8 dimensions)</i>			
17	dinuc.CG	-	CG dinucleotide frequency
18	dinuc.GC	-	GC dinucleotide frequency
19	dinuc.AT	-	AT dinucleotide frequency
20	dinuc.TA	-	TA dinucleotide frequency
21	dinuc.GG	-	GG dinucleotide frequency
22	dinuc.CC	-	CC dinucleotide frequency
23	dinuc.AA	-	AA dinucleotide frequency
24	dinuc.TT	-	TT dinucleotide frequency
<i>Type Indicator (1 dimension)</i>			
25	is_rna	0/1	RNA indicator (1) vs DNA (0)

Table S9: 320-dimensional fingerprint/embedding sources by material category.

Category	Method	Native Dim	Final Dim	Reference
Inorganic	MAGPIE composition	132	320 (zero-padded)	Ward et al., 2016
Molecules	Morgan fingerprint (r=2)	256	320 (zero-padded)	RDKit
Polymers	Morgan FP of repeat unit	256	320 (zero-padded)	RDKit
Biomolecules	ESM-2 mean pooling	320	320	Lin et al., 2023
DNA/RNA	DNABERT-S embedding	256	320 (zero-padded)	Zhou et al., 2024

- 2970
2971
2972
2973
2974
2975
2976
2977
2978
2979
2980
2981
2982
2983
2984
2985
2986
2987
2988
2989
2990
2991
1. **Electrical edges:** Connect source, drain, and substrate to the channel, representing carrier transport pathways.
 2. **Capacitive edges:** Model gate-channel coupling through the dielectric stack. The topology varies by device design:
 - Standard: gate \leftrightarrow dielectric \leftrightarrow channel
 - Remote: gate \leftrightarrow electrolyte \leftrightarrow channel
 - Floating gate: gate \leftrightarrow dielectric \leftrightarrow floating gate \leftrightarrow dielectric \leftrightarrow channel
 - Dual gate: parallel top and bottom gate pathways
 3. **Chemical edges:** Encode the sensing chain from channel through surface functionalization and probe to the detection target, plus target-medium and medium-channel interactions. We additionally include expert-recommended edges for probe-medium interactions (affecting probe stability and conformation), direct target-channel sensing (relevant for small molecules bypassing Debye screening), and probe-channel charge transfer.
 4. **Process edges:** Directed edges from annealing parameters to the channel, encoding thermal treatment effects on material properties.
 5. **Condition edges:** Directed edges from operating conditions to channel, gate, and test medium nodes.
 6. **Environment edges:** Connect electrolyte to test medium in remote-gate designs, modeling the shared solution environment between the gating electrolyte and the sensing medium.

2992 All chemical, electrical, capacitive, and environment edges are bidirectional to allow message passing in both directions, while process and condition edges are unidirectional (conditions affect components, not vice versa).

2996 C.2.3 PHYSICS-AWARE DATA AUGMENTATION

2997 Following the data protocol established in prior work [Ferreira et al. \(2025\)](#), we address the inherent size limitation of FET sensor datasets extracted from literature. To improve model generalization without introducing physically implausible samples, we develop a physics-aware data augmentation strategy that perturbs device parameters within experimentally reasonable bounds while enforcing domain constraints.

3003 **Label-safe Augmentation via Wide Classification Bins** A key insight enabling aggressive augmentation is that our classification bins span multiple orders of magnitude. For example, the lower detection limit (LDL) Class 1 spans from 10^{-6} to 10^0 ppm—a $10^6\times$ range. This wide binning provides a substantial safety margin: perturbations causing less than $10\times$ change in the target metric will not alter the class label. We exploit this property by designing perturbations that induce at most $\sim 20\text{--}25\%$ changes in device performance metrics, well within the label-safe regime.

3009 The physical basis for each perturbation’s impact can be estimated from device physics. For instance, gate dielectric thickness d affects transconductance as $g_m \propto C_{ox} \propto 1/d$, so a 25% thickness perturbation yields approximately 20% change in sensitivity—safely within the bin boundaries. Similarly, pH perturbations of $\pm 0.6\text{--}0.8$ units induce surface potential changes of $\Delta\psi \approx 59 \text{ mV} \times \Delta\text{pH} \approx 35\text{--}47 \text{ mV}$, which have negligible impact on detection limits for non-pH sensors.

3015 **Continuous Parameter Perturbation** For each original sample, we generate augmented variants by applying bounded perturbations to continuous parameters. Perturbation factors are sampled from truncated Gaussian distributions centered at unity (for multiplicative noise) or zero (for additive noise), with bounds reflecting typical experimental uncertainties (Table [S10](#)).

3020 **Physical Constraints** To ensure augmented samples remain physically realizable, we enforce constraints that reflect fundamental physical laws:

- 3022
3023
- **pH range:** Values are explicitly clipped to $[0, 14]$, corresponding to the thermodynamic limits of proton activity in aqueous solutions ($[\text{H}^+] = 1 \text{ M to } 10^{-14} \text{ M}$).

Table S10: Perturbation bounds for physics-aware data augmentation.

Parameter	Perturbation Type	Bounds	Physical Basis
Temperature (gas sensors)	Multiplicative	$\times 0.92\text{--}1.08$	$\pm 8\%$ calibration uncertainty
Temperature (biosensors)	Additive	$\pm 4^\circ\text{C}$	Physiological range variation
pH (biosensors)	Additive	± 0.6	Buffer preparation tolerance
pH (liquid sensors)	Additive	± 0.8	Environmental variation
Dielectric thickness	Multiplicative	$\times 0.85\text{--}1.15$	$\pm 15\%$ deposition variation
Substrate thickness	Multiplicative	$\times 0.85\text{--}1.15$	Wafer tolerance
Annealing temperature	Multiplicative	$\times 0.92\text{--}1.08$	Furnace calibration
Annealing time	Multiplicative	$\times 0.85\text{--}1.15$	Process window

- **Temperature and thickness:** These parameters use multiplicative perturbation factors ($\times 0.85\text{--}1.15$), which inherently preserve the sign of the original positive values, ensuring temperatures remain above absolute zero and layer dimensions remain positive.

These constraints are critical: without them, augmentation can generate samples that violate fundamental physical laws (e.g., negative pH, sub-zero temperatures), leading models to learn spurious correlations from impossible device configurations.

Discrete Augmentation Beyond continuous perturbations, we apply discrete transformations that exploit known physical symmetries and equivalences in FET sensor design:

- **Source/drain interchange:** In the absence of asymmetric doping or geometry, source and drain electrodes are physically interchangeable—the distinction is purely a measurement convention defining current direction. Our dataset confirms this symmetry: 98.5% of samples have identical source and drain materials.
- **Gate stack flip:** In dual-gate and floating-gate architectures, the capacitive coupling between gate stacks and the channel follows symmetric physics. Top and bottom gate/dielectric layer orderings can be reversed without altering the fundamental electrostatic control mechanism.
- **Inert atmosphere substitution:** For carrier gases (in gas sensors) and annealing atmospheres, chemically inert species (N_2 , Ar, He, Ne) serve identical purposes—preventing oxidation and providing a controlled environment. We replace nitrogen-based atmospheres with other inert gases using their actual material descriptors from PubChem.

These discrete augmentations generated 3,176 additional samples, each grounded in established semiconductor device physics rather than arbitrary permutations.

Material Descriptor Perturbation Material properties retrieved from databases carry inherent uncertainties from experimental measurements and computational predictions. To improve model robustness against these uncertainties, we perturb material representations:

- **Macroproperty noise:** Gaussian noise ($\sigma = 0.08$, i.e., $\pm 8\%$) is added to the 25-dimensional macroproperty vectors, reflecting typical measurement and prediction uncertainties in database-reported material properties.
- **Fingerprint perturbation:** For binary Morgan fingerprints, we apply stochastic bit flips with probability 0.08. For continuous embeddings (ESM-2, DNABERT-S), additive Gaussian noise is applied.

Baseline Augmentation Strategies To validate the importance of physics-aware constraints, we compare against two random augmentation baselines (Table S11):

- **Same-magnitude random:** Uses identical perturbation bounds as physics-aware augmentation but removes physical constraints (allowing $\text{pH} < 0$, negative temperatures) and discrete augmentations. This ablation isolates the contribution of domain constraints from perturbation magnitude.

Table S11: Comparison of augmentation strategies. Physics-aware augmentation enforces physical constraints and includes discrete transformations, while random baselines ablate these components.

Property	Physics	Same-mag	Destruct.
Perturbation mag.	$\pm 8\text{--}15\%$	$\pm 8\text{--}15\%$	$\pm 80\text{--}200\%$
Physical constraints	✓	×	×
Discrete augment.	✓ (3176)	×	×
Material noise	$\pm 8\%$	$\pm 8\%$	$\pm 50\%$
Variants/sample	2	4	6

- **Destructive random:** Applies large perturbations ($\pm 80\text{--}200\%$) without constraints, representing naive augmentation that ignores physical plausibility.

Empirically (5-fold CV on LDL/UDL/Sensitivity), all three augmentations achieve high accuracy on the augmented test splits (≥ 0.95), but only the physics-aware strategy transfers to held-out original data. Physics-aware augmentation yields Original accuracy of roughly 0.88/0.85/0.92 on the three tasks, whereas random and same-magnitude baselines collapse to $\sim 0.72/0.67/0.79$ and $\sim 0.73/0.67/0.79$, respectively. This gap demonstrates that domain constraints—not perturbation magnitude—are critical for preserving the real-device distribution.

Rationale. Models trained on physics-aware augmentation generalize to the held-out original distribution, while those trained on magnitude-matched random noise do not. This indicates the GNN is learning physically meaningful device–sensing relationships instead of overfitting to synthetic artifacts. Domain knowledge in the augmentation (valid ranges, symmetric transformations, inert substitutions) constrains message passing to realistic charge/chemical pathways and keeps fingerprint fields semantically aligned with real sensors. Unconstrained noise allows shortcut features that break on real devices. The strong Original-set gains therefore substantiate both the augmentation procedure and the downstream model’s use of domain priors.

Evaluation Protocol To validate that augmented data can serve as a legitimate training source, we adopt an evaluation protocol following prior work on spiking GNNs for FET sensors [Ferreira et al. \(2025\)](#). The augmented dataset is split 80:20 for training and testing, while the *entire original dataset* is reserved as a held-out test set.

The primary evaluation metric is **held-out accuracy**: classification accuracy on the complete original dataset. This metric measures whether a model trained on augmented data has learned patterns that transfer to real experimental samples. If the augmentation strategy preserves the underlying physical relationships, the model should generalize to original data; if augmentation introduces artifacts, the model will fail on real experiments.

Comparative experiments between physics-aware augmentation and random baselines (presented in Section [D](#)) demonstrate that physics-aware augmentation substantially outperforms random perturbation strategies on held-out accuracy. This validates that physics-aware augmented data captures genuine structure in the FET sensor design space and can be reliably used for training predictive models. Consequently, the held-out accuracy reported for downstream model benchmarking provides a meaningful measure of generalization to real-world sensing experiments.

C.2.4 GRAPH NEURAL NETWORK TRAINING

Graphs built in Section [C.2.2](#) (physics-informed augmented set as default) are used to train multi-class classifiers for three tasks: lower detection limit (LDL), upper detection limit (UDL), and sensitivity. Each graph carries both the “augmented” label (training distribution) and the original unperturbed label (held-out robustness evaluation).

Classification Task Definition Following prior work [Ferreira et al. \(2025\)](#), we discretize continuous performance metrics into three ordinal classes to enable robust classification despite measurement uncertainties inherent in literature-reported values. Class boundaries are defined on logarithmic scales to span multiple orders of magnitude (see Figure [S17](#)):

- **LDL (Lower Detection Limit):** Class 0 ($< 10^{-6}$ ppm, best), Class 1 (10^{-6} – 10^0 ppm), Class 2 ($> 10^0$ ppm, worst). Lower values indicate better sensitivity to trace analytes.
- **UDL (Upper Detection Limit):** Class 0 ($< 10^0$ ppm, worst), Class 1 (10^0 – 10^4 ppm), Class 2 ($> 10^4$ ppm, best). Higher values indicate broader dynamic range.
- **Sensitivity:** Class 0 (low), Class 1 (medium), Class 2 (high, best). Boundaries depend on the sensitivity definition (e.g., $\Delta R/R_0$, $\Delta I/I_0$) and are normalized per metric type.

This wide binning provides label stability under data augmentation: perturbations causing $< 10\times$ metric changes remain within the same class, enabling physics-aware augmentation without label corruption.

Model Architecture Our primary model is a residual heterogeneous GNN with two branches:

- **GNN branch:** For each edge type (s, r, d) in the heterogeneous graph, we use relation-specific convolution operators aggregated across relations (Eq. 4). Self-loops are optionally added at the graph level; we disable internal self-loop insertion for GATv2Conv. We apply one hetero-convolution layer Schlichtkrull et al. (2018a) with jump-knowledge residuals Xu et al. (2018) and GCNII-style Chen et al. (2020) re-parameterization (Eq. 5), followed by LayerNorm, ReLU, and an attention-based readout. Node inputs are 51-dimensional vectors composed of 25 macroproperties + 25 missing-data masks + 1 material-count scalar (Section C.2.1); the 320-dimensional fingerprints are *not* fed here.
- **Fingerprint branch:** A small MLP that processes the 320-dimensional fingerprints for five fields most relevant to the chemical sensing interface: channel (transducer), detect_target (analyte), probe_material (recognition element), test_medium (sample matrix), and surface_functionalization (interface modifier). Its output is fused with the GNN logits through a learnable gate, initialized to a fixed value and updated in the second training stage.

This architecture preserves physical message-passing paths while allowing complementary fingerprint evidence to modulate the prediction.

Training Protocol We adopt a two-stage schedule: (1) pre-train the GNN branch for E_{gnn} epochs; (2) freeze the GNN and train the fingerprint gate/MLP for E_{fp} epochs. For each fold we report four metrics (accuracy, macro-F1, precision, recall) on (a) augmented train, (b) augmented test, and (c) original held-out graphs; final results are mean \pm std across 5-fold stratified cross-validation.

Key hyperparameters (hidden dimension, dropout, learning rate, weight decay, gate initialization, E_{gnn} , E_{fp}) are selected via Bayesian optimization (BO) with adaptive level-set estimation Zhang et al. (2023), which filters for a high-confidence region of interest (ROI) as a superlevel-set of a Gaussian process surrogate, around a manually tuned seed configuration. The anchor used for all ablations is: hidden_dim=128, dropout=0.2, lr= 1×10^{-4} , weight_decay= 5×10^{-5} , GCNII $\alpha = 0.15$, pretrain epochs $E_{\text{gnn}} = 40$, FP epochs $E_{\text{fp}} = 20$, gate_init=0.25, self-loops on. All ablation variants reuse this configuration for fairness. Batches are class-stratified. Optimization uses AdamW Loshchilov & Hutter (2019) with weighted cross-entropy loss (Eq. 7), implemented via nn.CrossEntropyLoss which includes the softmax internally. Unless otherwise stated, batch size is 32.

Controls and Ablations To attribute gains, we conduct 14 controlled ablation variants organized into six categories:

- **Architecture components:** toggling jump-knowledge (JK) connections and attention-based readout—full model (JK+Attention), w/o JK, w/o Attention, w/o both.
- **Branch isolation:** GNN-only (no fingerprint branch) and fingerprint-only (FP-only, no graph convolution) to quantify each branch’s contribution.
- **Alternative FP encoders:** replacing the default MLP fingerprint branch with Transformer or Spiking Neural Network (SNN) encoders, sharing anchor hyperparameters.
- **GNN depth:** varying the number of hetero-convolution layers (0, 1 [default], 2) to assess depth effects.

- 3186 • **GCNII residual strength:** varying the initial residual coefficient $\alpha \in \{0.05, 0.15, 0.25\}$;
- 3187 $\alpha = 0.15$ is the anchor default.
- 3188 • **Training variants:** toggling self-loops in message passing, and enabling inverse-frequency
- 3189 class weights for imbalanced labels.
- 3190

3191 Additionally, we compare against 17 tabular baselines trained on flattened descriptors, each tuned
 3192 via the same ROI-based BO protocol: Random Forest (RF), Extremely Randomized Trees (Ex-
 3193 traTrees), Decision Tree (DT), XGBoost, LightGBM, Gradient Boosting Decision Tree (GBDT),
 3194 CatBoost, AdaBoost, Logistic Regression (LogReg), Linear Support Vector Machine (LIN_SVM),
 3195 Radial Basis Function SVM (RBF_SVM), K-Nearest Neighbors (KNN), Multi-Layer Perceptron
 3196 (MLP), Linear Discriminant Analysis (LDA), Quadratic Discriminant Analysis (QDA), Ridge Clas-
 3197 sifier, and Gaussian Naive Bayes (GNB). We also include a standalone Spiking GNN (SGNN) base-
 3198 line from prior work [Ferreira et al. \(2025\)](#).

3199 Unless otherwise stated, “anchor” refers to the best BO-selected residual GNN ($\alpha = 0.15$,
 3200 JK+Attention, 1 GNN layer, self-loops enabled) with key fingerprints and physics-informed aug-
 3201 mentation; all ablation variants reuse its hyperparameters for fairness.

3202

3203 C.3 METHOD AND TECHNICAL DETAILS FOR “TRANSLATION”

3204

3205 C.3.1 VIRTUAL SCREENING METHODOLOGY

3206 **Candidate Molecule Library** The screening library comprised 123,239,643 molecules obtained
 3207 from the PubChem Compound database (CID-SMILES file, accessed December 2025). Each
 3208 molecule was represented by:

3209

- 3210 • A 25-dimensional molecular descriptor vector containing physicochemical properties in-
 3211 cluding molecular weight, topological polar surface area (TPSA), hydrogen bond donor/ac-
 3212 ceptor counts, rotatable bond count, XLogP, and 3D steric parameters from PubChem.
- 3213 • A 256-bit Morgan fingerprint (radius 2) generated using RDKit for structural encoding.

3214

3215 The library was partitioned into 62 chunks of 2,000,000 molecules each, further split into parts of
 3216 500,000 molecules for parallel processing. A total of 235 parallel jobs were executed on GPU nodes
 3217 (NVIDIA A40), achieving a throughput of approximately 40 molecules per second per job.

3218

3219 **Graph Construction for Inference** For each candidate molecule, we constructed heterogeneous
 3220 graphs following the DTE-GNN architecture. The template graph was derived from the training data
 3221 corresponding to the PFOS detection experiments in [Wang et al. \(2025\)](#), containing the following
 3222 node types:

- 3223 • **Probe material:** Replaced with each candidate molecule’s features.
- 3224 • **Detect target:** Set to PFOS, PFOA, dodecylsulfonic acid (DDS), or trichloroacetic acid
 3225 (TCAA) depending on the screening task.
- 3226 • **Substrate:** Fixed to the experimental configuration (kept from template).
- 3227 • **Condition:** Fixed measurement conditions including concentration ranges and environ-
 3228 mental parameters.

3229

3230 The edge structure connecting these nodes was preserved from the training template, representing
 3231 the physical relationships in the FET sensor system.

3232

3233 **Target Molecules for Selectivity Analysis** Four detect target molecules are used for comprehen-
 3234 sive selectivity screening (Table [S12](#)):

3235

3236 **Model Inference** Three pre-trained DTE-GNN models are applied to each candidate-target pair:

3237

- 3238 1. **LDL Model:** Predicts lower detection limit class (0: low/good, 1: medium, 2: high/poor).
- 3239 2. **UDL Model:** Predicts upper detection limit class (0: low/poor, 1: medium, 2: high/good).

Target	CID	Role	Description
PFOS	74483	Primary	Perfluorooctanesulfonic acid
PFOA	9554	Target	Perfluorooctanoic acid
DDS	3423265	Interferent	Dodecylsulfonic acid
TCAA	6421	Interferent	Trichloroacetic acid

Table S12: Target molecules used in selectivity screening. Neutral forms are used for GNN screening; anionic forms for DFT validation (Appendix C.3.2).

3. **Sensitivity Model:** Predicts sensitivity class (0: low/poor, 1: medium, 2: high/good), evaluated for both percentage-based and mV-based sensitivity metrics.

Each model outputs a probability distribution $[p_0, p_1, p_2]$ over three classes. Inference is performed in batches of 2,048 molecules using PyTorch and PyTorch Geometric on CUDA-enabled GPUs.

Scoring Functions

SINGLE-TARGET SCORE For a given probe-target pair, the composite score reflects the joint probability of achieving desirable performance across all three metrics:

$$S_{\text{target}} = P(\text{LDL} = 0) \times P(\text{UDL} = 2) \times P(\text{Sensitivity} = 2) \quad (8)$$

where class 0 for LDL indicates low detection limit (desirable), and class 2 for UDL and sensitivity indicates high dynamic range and high sensitivity, respectively.

SELECTIVITY SCORE To identify probes with high specificity for PFAS compounds over common interferents, we computed the selectivity score as:

$$S_{\text{selectivity}} = \frac{S_{\text{PFOS}} \times S_{\text{PFOA}}}{S_{\text{DDS}} \times S_{\text{TCAA}}} \quad (9)$$

This formulation rewards candidates that:

- Achieve high scores for both PFOS and PFOA (numerator, representing target PFAS response).
- Achieve low scores for both DDS and TCAA (denominator, representing interferent response).

Candidates were filtered to require a minimum target score of 0.01 and selectivity ratio greater than 1.0 before ranking.

Computational Resources The complete screening of 123M molecules across four targets required:

- 940 GPU-hours (235 jobs \times 4 targets \times approximately 1 hour each).
- Storage: 120 GB for result files (30 GB per target).
- Peak memory: 4-8 GB per job during chunk processing.

Result Aggregation Results were aggregated using a memory-efficient streaming algorithm:

1. Process one chunk at a time, loading corresponding result files from all four targets.
2. Compute selectivity scores for molecules present in all four result sets.
3. Maintain a min-heap of size 1,000 to track top candidates without storing all results in memory.
4. Output final rankings with complete prediction details for downstream analysis.

C.3.2 DFT COMPUTATIONAL METHODOLOGY FOR HOST-GUEST BINDING ENERGY CALCULATIONS

Scope and Rationale The DFT calculations in this work serve as qualitative validation of the DTE-GNN screening predictions, specifically to compare binding selectivity between the model-predicted probe candidates and the experimentally validated β -Cyclodextrin baseline [Wang et al. \(2025\)](#). Given the broad scope of the T³ framework spanning text mining, GNN modeling, and virtual screening across millions of candidates, we adopt a computationally efficient grid-based conformational sampling approach with implicit solvation rather than expensive classical molecular dynamics (MD) simulations with explicit solvent. This choice is justified for our comparative screening purpose: systematic grid sampling adequately captures the dominant binding modes for rigid macrocyclic hosts, while the SMD implicit solvation model [Marenich et al. \(2009\)](#) provides reliable solvation free energies at a fraction of the computational cost. More rigorous MD-based free energy calculations could be pursued in future studies for the most promising candidates identified here.

Initial Structure Generation Host-guest inclusion complex configurations are systematically generated through geometric sampling to ensure comprehensive coverage of the binding conformational space. For each host-guest pair, three types of configurations are constructed:

1. **Vertical insertion configurations:** The guest molecule is aligned along the host cavity axis and positioned at nine different insertion depths (-8, -6, -4, -2, 0, +2, +4, +6, +8 Å relative to the cavity center). At each depth, two orientations are sampled by flipping the guest molecule 180 degrees, resulting in 18 configurations per pair.
2. **Surface-lying configurations:** The guest molecule is placed horizontally on both the upper and lower surfaces of the host cavity. At each surface, four rotational orientations (0, 90, 180, 270 degrees around the cavity axis) are sampled, yielding 8 configurations per pair.
3. **Side-binding configurations:** The guest molecule is positioned approaching the host from the side at six azimuthal angles (0, 60, 120, 180, 240, 300 degrees), generating 6 configurations per pair.

The host cavity axis is determined using principal component analysis (PCA) on the ring atoms defining the macrocyclic cavity. The direction corresponding to the smallest eigenvalue is identified as the cavity normal vector. Similarly, the guest molecule principal axis is determined as the direction of largest structural variance via PCA.

All initial structures are placed in a cubic simulation box of 50 Å side length to ensure adequate separation from periodic images.

Geometry Optimization A multi-level optimization protocol is employed to efficiently explore the potential energy surface:

1. **Semi-empirical pre-optimization:** Initial structures are first optimized using the PM6 semi-empirical method to remove atomic clashes and obtain reasonable starting geometries.
2. **DFT geometry optimization:** The PM6-optimized structures are further refined using density functional theory at the B3LYP/6-31G(d) level with the SMD implicit solvation model (solvent = water). This step provides accurate equilibrium geometries while maintaining computational tractability for the large host-guest systems.

Single-Point Energy Calculations Final electronic energies are computed at a higher level of theory using single-point calculations on the B3LYP-optimized geometries. The M06-2X functional with the 6-31+G(d) basis set is employed, combined with the SMD solvation model [Marenich et al. \(2009\)](#) to account for aqueous solvation effects. This computational protocol follows established approaches for PFOS thermodynamic calculations [Montero-Campillo et al. \(2010\)](#); [Giroday et al. \(2014\)](#). The M06-2X functional is selected for its reliable description of non-covalent interactions, including dispersion and hydrogen bonding, which are crucial for accurate host-guest binding energetics.

Configuration Selection and Binding Energy Calculation For each host-guest pair, the configuration with the lowest total electronic energy from the single-point calculations is selected as the representative binding geometry. The binding energy is calculated as:

$$\Delta E_{\text{bind}} = E_{\text{complex}} - E_{\text{host}} - E_{\text{guest}} \quad (10)$$

where E_{complex} , E_{host} , and E_{guest} are the single-point energies of the host-guest complex, isolated host molecule, and isolated guest molecule, respectively. All energies are converted from Hartree to kcal/mol using the conversion factor 627.509474 kcal/mol per Hartree.

Treatment of Protonation States To investigate the influence of guest molecule protonation state on binding affinity, calculations are performed for both neutral guest molecules (protonated forms: PFOA, PFOS, DDS-H, TCAA) and their corresponding deprotonated anionic forms (PFOA⁻, PFOS⁻, DDS⁻, TCA⁻). For anionic systems, geometry optimization and single-point calculations are performed with a total charge of -1 .

Software and Computational Resources All quantum chemical calculations are performed using Gaussian 16. Geometry optimizations employ default convergence criteria. Single-point calculations use tight SCF convergence. Structure generation and analysis scripts are implemented in Python using NumPy and SciPy libraries.

Electrostatic Potential Surface Analysis To visualize and quantify the electrostatic environment of host-guest complexes, electrostatic potential (ESP) surfaces are computed using the Gaussian 16 `cubegen` utility. For each system, a dedicated single-point calculation is performed at the ω B97XD/6-311+G(2d,2p) level with SMD implicit solvation (solvent = water) on the lowest-energy binding configuration, with `density=current` and `pop=full` keywords to ensure the ESP is evaluated from the converged SCF density. The SCF electron density and electrostatic potential are then evaluated on a uniform $100 \times 100 \times 100$ grid via `cubegen`. The van der Waals surface is defined as the $\rho = 0.001$ a.u. electron density isosurface, following the convention of Murray & Politzer (2011). ESP values on this isosurface are extracted and converted from atomic units (Hartree/e) to electronvolts (1 Hartree = 27.2114 eV). Quantitative descriptors include the mean surface ESP (\bar{V}), standard deviation (σ), fraction of positive/negative surface area, and the Politzer electrostatic balance index Murray & Politzer (2011). ESP surfaces are rendered using PyMOL with a color scale of -0.05 to $+0.05$ a.u. (red-white-blue).

D AUGMENTATION STRATEGY ABLATION

Following the data protocol from Ferreira et al. (2025), we validate the effectiveness of physics-aware data augmentation by comparing against two random baselines: same-magnitude random (identical perturbation bounds but without physical constraints or discrete augmentations) and destructive random (large perturbations of ± 80 – 200% without constraints). All models are trained on augmented data and evaluated on the held-out original dataset to assess generalization to real experimental samples.

Figure S25 summarizes the results across LDL, UDL, and Sensitivity prediction tasks. Physics-aware augmentation consistently outperforms both random baselines by substantial margins. On basic metrics (Accuracy, F1, Precision, Recall), physics-aware augmentation achieves 13–18% absolute improvements over the best random baseline across all three tasks, with UDL showing the largest gap (+17.7% accuracy, +17.9% F1). Notably, same-magnitude random performs comparably to destructive random despite using much smaller perturbations, demonstrating that the performance gap is not simply due to noise magnitude. Rather, physics-aware augmentation generates an entire dataset of physically plausible FET sensor configurations—each augmented sample represents a device that could realistically exist—whereas random perturbations produce nonsensical combinations (e.g., negative pH, impossible material properties) that corrupt the learned representations. This validates that our augmentation strategy creates meaningful synthetic data grounded in semiconductor device physics, enabling models to learn transferable patterns from an expanded but physically coherent training distribution.

3402
3403
3404
3405
3406
3407
3408
3409
3410
3411
3412
3413
3414
3415
3416
3417
3418
3419
3420
3421
3422
3423
3424
3425
3426
3427
3428
3429
3430
3431
3432
3433
3434
3435
3436
3437
3438
3439
3440
3441
3442
3443
3444
3445
3446
3447
3448
3449
3450
3451
3452
3453
3454
3455

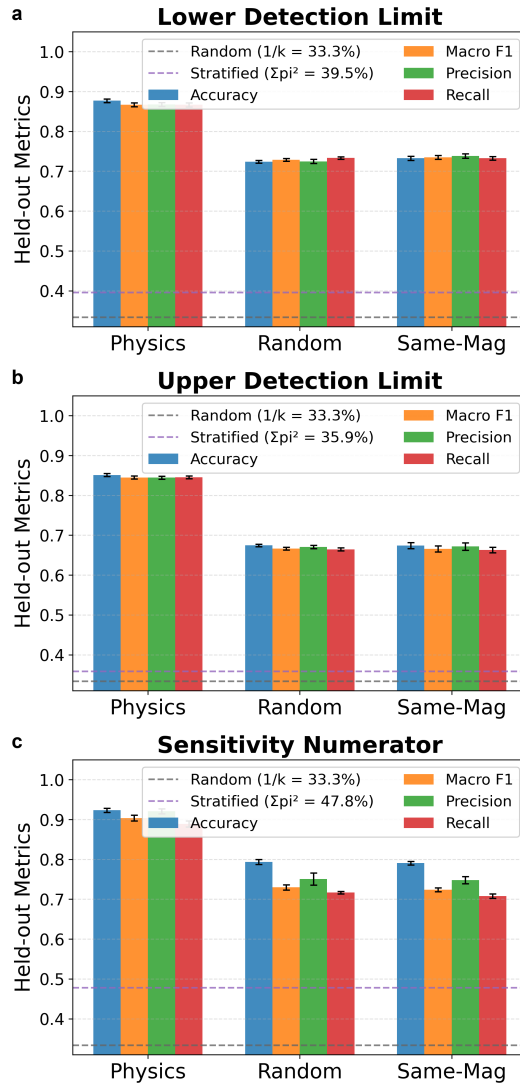


Figure S25: Augmentation strategy comparison on held-out original data. Physics-aware augmentation vs. same-magnitude random and destructive random baselines across LDL, UDL, and Sensitivity prediction tasks. Metrics shown: Accuracy, macro-F1, Precision, and Recall.

E DTE-GNN ARCHITECTURE ABLATION

To understand the contribution of each architectural component in DTE-GNN, we conduct systematic ablation studies across 14 configurations varying readout mechanisms, branch composition, fingerprint encoders, network depth, GCNII residual strength (α), and training strategies (Figure 4b–c).

Multi-Task Ranking Methodology. Since different configurations may excel on different prediction tasks, we adopt a cross-task ranking approach following established practices for comparing classifiers over multiple datasets Demsar (2006). For each configuration c and task t , we compute an optimistic score $s_{c,t} = \mu_{c,t} + \sigma_{c,t}$, where $\mu_{c,t}$ is the mean accuracy and $\sigma_{c,t}$ is the standard deviation across cross-validation folds. This UCB-style scoring, inspired by Gaussian process optimization Derrac et al. (2011), rewards configurations that achieve high mean performance while accounting for estimation uncertainty—an “optimism in the face of uncertainty” principle that balances exploration and exploitation. We rank configurations within each task by their optimistic scores and compute the average rank across all three tasks (LDL, UDL, Sensitivity). The configuration with the lowest average rank is considered the most consistently effective. Under this methodology, the full DTE-GNN model (with JK connections, attention pooling, and $\alpha = 0.15$) achieves the best cross-task ranking in both alpha ablation and component ablation studies, validating it as the most robust architecture choice (Figure S19).

The most striking finding is that the GNN branch is essential: removing it entirely (fingerprint (FP) Only) causes performance to collapse to near-random levels (31–42% accuracy vs. 87–92% for the full model), confirming that the heterogeneous graph structure captures the core device-performance relationships. In contrast, the fingerprint branch provides modest but consistent improvements; GNN Only configurations perform within 1% of the full model, indicating that macroscopic physical properties encoded in graph nodes carry the dominant predictive signal.

Regarding readout mechanisms, the full model with both jumping-knowledge (JK) connections and attention-based pooling achieves the most consistent performance across all tasks. While individual ablations may show marginal improvements on specific tasks (e.g., removing JK yields 88.2% vs. 87.7% accuracy on LDL alone), cross-task evaluation reveals that such gains do not generalize—the full model maintains the best average ranking when evaluated across LDL, UDL, and Sensitivity simultaneously. Network depth matters: zero GNN layers (Depth: 0) substantially degrades performance (5–7% accuracy drop), while two layers perform comparably to one layer, indicating that a single message-passing step suffices to propagate information across the physics-informed graph topology. The GCNII residual strength ($\alpha \in \{0.05, 0.15, 0.25\}$) and fingerprint encoder choice (MLP vs. Transformer vs. SNN) show minimal impact (<0.5% variation), suggesting the architecture is robust to these hyperparameters. Finally, class weighting and self-loop removal provide no consistent benefit, validating our default training configuration.

F PROMPT EVOLUTION

This appendix presents the human-authored initial prompts alongside their autonomously optimized counterparts for all five field categories: the core fields (8 fields) and four extended field groups. For the four extended field groups, we embed the representative top-performing prompts (M8, inference=Deepseek-70B, critique=GPT-oss-120B). Across these prompts, the optimizer generally learned to: (i) tightly scope the extraction window to the main article while ignoring cited abstracts/supplementary fluff; (ii) normalize units and apply explicit defaults (thickness/time/concentration) instead of heuristic guessing; (iii) disambiguate primary device elements from auxiliary layers (e.g., gate dielectric vs. encapsulation vs. reference electrode); (iv) collapse multiple material or device variants into a single representative record when the paper only tweaks minor formulations; (v) extract numeric values directly with required conversions (e.g., pH→[H⁺], ppm from P/P_v, seconds from minutes); and (vi) forbid inference—fields remain "" unless explicitly stated. These behaviors collectively reduce hallucination, enforce consistency, and preserve the most representative device configuration per paper.

Concrete improvements observed across the five prompt families include: (1) *Scope control*—original prompts often swept in cited abstracts or supplementary blurbs; optimized prompts now

3510 restrict extraction to the main article body (and, when needed, a specified table or subsection). (2)
 3511 *Unit rigor*—where initial prompts left units free-form, optimized prompts demand explicit conver-
 3512 sions (e.g., minutes→seconds for response/recovery times; P/P_v→ppm; pH→[H⁺] for bounds) or
 3513 defaults when absent. (3) *Device disentangling*—early instructions blurred substrate/dielectric/en-
 3514 capsulation; optimized prompts explicitly separate gate dielectric from protective coatings and re-
 3515 quire empty strings when not stated. (4) *Variant consolidation*—initial versions invited multiple
 3516 records for minor material tweaks; optimized prompts collapse variants into a single representa-
 3517 tive record unless the paper truly reports distinct devices. (5) *Field-specific cues*—material-layer
 3518 prompts add dopant handling and thickness defaults; sensitivity/response prompts require table pars-
 3519 ing and unit rounding; synthesis/thermal prompts stress atmosphere/time/temperature capture while
 3520 ignoring non-functional processing layers. Collectively, these edits yield cleaner, reproducible JSON
 3521 with fewer missing-unit errors and no fabricated values.

3522 F.1 CASE STUDY: PROMPT EVOLUTION TRAJECTORY

3524 To illustrate how TextGrad-based optimization progressively refines extraction prompts, we trace a
 3525 single representative run (M8, DeepSeek-70B inference, GPT-oss-120B critique) across 20 training
 3526 rounds on the core fields. Under our tournament-based selection mechanism (Section C.1), the
 3527 optimized candidate replaces the current prompt *only* when it achieves a higher committee score
 3528 on the same minibatch; otherwise the existing prompt is retained unchanged. In this run, 8 of
 3529 20 proposals were accepted (40%), reflecting the conservative selection pressure that filters out
 3530 regressions while accumulating genuine improvements.

3532 **Stage 1: General extraction hygiene (Rounds 2–6).** The earliest accepted changes address
 3533 generic extraction errors that manifest across all paper types.

- 3535 • **Round 2 — Verbatim probe-material names.** The human prompt specifies "(probe
 3536 chemical name)", which led the LLM to paraphrase or abbreviate materials (e.g.,
 3537 producing "zinc oxide" instead of the paper's "ZnO nanowires"). The optimizer re-
 3538 placed this with "(exact probe material name as it appears in the
 3539 paper, copy verbatim)", eliminating name mismatches. This instruction persisted
 3540 unmodified through all 18 subsequent rounds.
- 3541 • **Round 4 — Unit-bias removal for detection limits.** The original placeholder "xx ppm
 3542 (or corresponding units)" acted as an implicit prior, biasing the LLM toward
 3543 converting all concentrations to ppm. The optimizer replaced it with "value with its
 3544 original unit (or leave empty)", preserving the diversity of units (nM, ppb,
 3545 vol%, etc.) reported across the sensor literature.
- 3546 • **Round 6 — Scope control against reference hallucination.** A single parenthetical
 3547 clause was added to the opening instruction: "(excluding bibliography and
 3548 reference list)". This prevented the LLM from extracting sensor specifications
 3549 mentioned in cited references rather than the paper's own experiments—a source of phan-
 3550 tom records. This round produced the largest single-round committee-score improvement
 3551 (+0.117).

3552 **Stage 2: Domain-specific disambiguation (Rounds 7–12).** Once general extraction hygiene is
 3553 established, subsequent improvements target sensor-science-specific ambiguities that arise only for
 3554 particular paper types.

- 3556 • **Round 7 — Humidity sensor classification.** The sensor_type field was augmented with
 3557 the rule: "water vapour in air counts as gas". This resolved systematic
 3558 misclassification of humidity sensors, which were previously labeled as "liquid" due to the
 3559 presence of water.
- 3560 • **Round 12 — Detection-limit fallback from tested concentrations.** Many papers omit
 3561 an explicit limit of detection but report a tested concentration range. The optimizer
 3562 added: "if not directly stated, use the smallest concentration
 3563 tested". This reduced missing-value rates for detection limits without introducing fab-
 ricated numbers.

Between these breakthroughs, two stabilization plateaus occurred (Rounds 8–11 and 13–16), during which proposals were consistently rejected. Notably, several rejected changes directly contradicted previously learned rules—for instance, Round 13 proposed expanding chemical abbreviations to full IUPAC names, conflicting with the “copy verbatim” rule from Round 2. The tournament mechanism correctly preserved the earlier beneficial modification.

Stage 3: Numerical conversion rules (Rounds 17–20). The most domain-specific improvements emerged in the final rounds, when rarer paper types (e.g., pH sensors) exposed new error classes.

- **Round 17 — pH-to-concentration conversion formula.** pH sensors report detection ranges in pH units (e.g., “pH 2–12”), whereas the ground-truth database records molar concentrations. The optimizer discovered and embedded the conversion rule: "calculate lower_detection_limit = $10^{-\text{upper_pH}}$ M", including the non-trivial inversion—upper pH maps to *lower* concentration—without any explicit chemical instruction. This illustrates TextGrad’s capacity to learn domain-specific numerical transformations purely from error feedback between LLM outputs and expert annotations.
- **Round 20 — Temperature missing-data handling.** A final instruction was added to output an empty field rather than a default numeric value when no operating temperature is stated, correcting the LLM’s tendency to hallucinate “25 °C” as a placeholder.

Table S13 summarizes the trajectory. The evolution proceeds from generic extraction hygiene (verbatim copying, unit preservation, scope control) through domain-specific disambiguation (sensor classification, fallback logic) to numerical conversion rules (pH→concentration). Each accepted modification is strictly additive—no accepted change was ever reverted in a later round. The conservative tournament selection (60% rejection rate) plays a critical role: it prevents regression while allowing incremental accumulation of domain knowledge that would be difficult to engineer manually. Table S14 summarizes the 28 extraction fields organized into five prompt categories.

Table S13: Prompt evolution trajectory for M8 on core fields (20 rounds). Accepted changes are grouped by the type of improvement.

Round	Improvement	Category	ΔScore
2	Verbatim material names	General	+0.000
4	Unit-bias removal	General	+0.012
5	Test medium disambiguation	General	+0.008
6	Scope: exclude bibliography	General	+0.117
7	Humidity sensor rule	Domain	+0.015
12	Detection-limit fallback	Domain	+0.053
17	pH→concentration formula	Numerical	+0.010
20	Temperature missing-data	Numerical	+0.023
<i>Rejected proposals</i>			12/20

F.2 INITIAL AND OPTIMIZED PROMPTS

CORE FIELDS (8 FIELDS)

Input: Full text of a Paper

Expected Output Example (by human expert):

```
{
  "records": [
    {
      "sensor_type": "gas",
      "detect_target": "ammonia",
      "lower_detection_limit": "4 ppm",
      "upper_detection_limit": "60 ppm",
      "probe_material": "6PTTP6(5,5'-Bis(4-hexylphenyl)-2,2'-bithiophene)
/8-3-NTCDI(N,N'-Bis(3-(perfluorooctyl)propyl)-1,4,5,8-
naphthalenetetracarboxylic diimide)",
```

3618 Table S14: Summary of LLM extraction fields by prompt category. A total of 28 fields are extracted
 3619 across five prompt families.
 3620

3621 Category	3622 Fields	3623 Field Names
3624 Core	3625 8	3626 sensor_type, detect_target, lower_detection_limit, upper_detection_limit, 3627 probe_material, test_operating_temperature, pH_value, test_medium
3628 Electrode Architecture	3629 4	3630 gate, source, drain, structure_design_type
3631 Material Layer	3632 7	3633 substrate, substrate.thickness, channel, dielectric_layer, di- 3634 electric_layer_thickness, surface_functionalization, struc- 3635 ture_dimensionality
3636 Sensitivity/Response	3637 4	3638 response_time, recovery_time, sensitivity_numerator, sensitiv- 3639 ity_denominator
3640 Synthesis/Thermal	3641 5	3642 annealing_temperature, annealing_time, annealing_atmosphere, hy- 3643 drothermal_temperature, hydrothermal_time
3644 Total	3645 28	

```

3636     "test_operating_temperature (celcius)": "25 \textdegree C",
3637     "pH_value": "-1",
3638     "test_medium": "air"
3639   }
3640 ]
3641 }

```

3642 Human initial prompt (4 fields)

3643 After reading the scientific publication full text provided above, please
 3644 try to generate formatted JSON file for extraction of important
 3645 information. Leave "" for not available fields.
 3646

3647 For instance:

```

3648 {
3649   "records": [
3650     {
3651       "sensor_type": "(gas/bio/liquid; please choose one)",
3652       "detect_target": "(target chemical name)",
3653       "lower_detection_limit": "xx ppm (or corresponding units)",
3654       "upper_detection_limit": "xx ppm (or corresponding units)",
3655       "probe_material": "(probe chemical name)",
3656       "test_operating_temperature (celcius)": "xx \textdegree C",
3657       "pH_value": "(use -1 for gas, otherwise a number, if pH sensor, use  

3658         range like xx-yy)",
3659       "test_medium": "(the name of the medium, e.g. air, water, or other  

3660         medium)"
3661     }
3662     (continue if the publication has recorded multiple different target  

3663     been detected)
3664   ]
3665 }

```

3663 Optimized prompt (M8, GPT-oss-120B critique, DS70B inference)

3665 After reading the scientific publication full text provided above, please
 3666 generate a formatted JSON file for extraction of important
 3667 information. Only use the main article text above; ignore any cited-
 3668 paper abstracts, supplementary sections, or other unrelated excerpts.
 3669 If the article describes multiple sensors, identify the sensor with
 3670 the smallest lower_detection_limit and create a single JSON object
 3671 for that sensor only. Do not create separate records for each
 material variant; combine them into one record representing the
 overall sensor system (use a generic class name for the probe

```

3672     material). When several material variants are tested for the same
3673     sensor, combine them into one record and describe the probe material
3674     using its generic class name (e.g., "diketopyrrolopyrrole fluorene
3675     copolymer"). Leave "" for not available fields.
3676
3677 For instance:
3678 {
3679   "records": [
3680     {
3681       "sensor_type": "(gas/bio/liquid; please choose one)",
3682       "detect_target": "(target chemical name)",
3683       "lower_detection_limit": "Extract the reported concentration value
3684         for the target chemical as written; provide only the numeric
3685         value with its unit, stripping any qualitative qualifiers such
3686         as 'below', 'above', or '~'.
3687         Prioritize any explicit detection-limit statement (e.g., 'the
3688         lowest detection limit is X ppm').
3689         If no explicit detection limit is given, use the smallest
3690         concentration that was experimentally tested (or leave '').
3691         If the paper provides a concentration range (e.g., 'from X ppm to Y
3692         ppm'), use the lower bound (X) as the lower detection limit.
3693         If the paper reports a 'limit of detection (LOD)', treat that value
3694         as the lower detection limit.
3695         If the value is given as a P/Pv ratio, convert it to ppm by
3696         multiplying the ratio by the analyte's equilibrium vapor
3697         pressure (Pv, kPa) and then by 1000 ppm / 101.3 kPa (~9.87).
3698         Output the numeric ppm value with the unit ppm.
3699         Otherwise preserve the exact numeric value and unit; do not convert
3700         , round, or infer other units.
3701         If the article gives a pH response range, convert the highest pH in
3702         that range to molar [H+] using [H+]=10^{-pH} M and record that
3703         value as the lower detection limit.",
3704       "upper_detection_limit": "The highest concentration of the target
3705         analyte that is explicitly reported or tested in the paper (i.e
3706         ., the maximum concentration examined in the experiments).
3707         Provide only the numeric value with its unit, stripping any
3708         qualitative qualifiers such as 'below', 'above', or '~'.
3709         Use the largest reported concentration (or leave '').
3710         If the paper provides a concentration range (e.g., 'from X ppm to Y
3711         ppm'), use the upper bound (Y) as the upper detection limit.
3712         If only a single value is given, leave 'upper_detection_limit'
3713         empty.
3714         If the value is given as a P/Pv ratio, convert it to ppm by
3715         multiplying the ratio by the analyte's equilibrium vapor
3716         pressure (Pv, kPa) and then by 1000 ppm / 101.3 kPa (~9.87).
3717         Output the numeric ppm value with the unit ppm.
3718         Otherwise preserve the exact numeric value and unit; do not convert
3719         , round, or infer other units.
3720         If the article gives a pH response range, convert the lowest pH in
3721         that range to molar [H+] using [H+]=10^{-pH} M and record that
3722         value as the upper detection limit.",
3723       "probe_material": "(material that directly interacts with the
3724         target gas, e.g., catalytic or gate metal layer; if no separate
3725         layer is described, use the gate-dielectric material; if
3726         several variants are reported, give the generic class name;
3727         otherwise leave '' if not available)",
3728       "test_operating_temperature (celcius)": "xx \textdegree C",
3729       // If the paper mentions an operating temperature (e.g.,
3730       // 'measurements were performed at 25 °C'), extract the numeric
3731       // value (with unit) and place it in 'test_operating_temperature
3732       // (celcius)'; if none, leave the field blank.
3733       "pH_value": "(if not directly given but can be reasonably inferred
3734         from the test medium, provide that inferred value; otherwise
3735         use -1; for gas sensors use -1; otherwise a numeric value; if
3736         pH sensor, use range like xx-yy)",

```

```

3726     "test_medium": "(the name of the medium, e.g., air, water, or other
3727         medium; if multiple media are mentioned, choose the one used
3728         for the primary sensor performance experiments, ignoring
3729         ancillary tests unless they are the only measurements reported)
3730     "
3731     // If multiple ranges are given, use the primary detection window:
3732     // from the lowest measurable concentration to the highest before
3733     // saturation.
3734 }
3735 (continue if the publication has recorded multiple different target
3736     been detected)
3737 ]
3738 }

```

EXTENDED FIELDS

Electrode Architecture (4 fields) Human initial prompt

After reading the scientific publication full text provided above, please try to generate formatted JSON file for extraction of important information. Leave "" for not available fields.

For instance:

```

3746 {
3747   "records": [
3748     {
3749       "gate": "(gate electrode material name)",
3750       "source": "(source electrode material name)",
3751       "drain": "(drain electrode material name)",
3752       "structure_design_type": "(e.g. planar, vertical, coplanar, back-
3753         gated, top-gated, interdigitated)"
3754     }
3755   ]
3756 }

```

(continue if the publication has recorded multiple different device architectures)

Optimized prompt (M8, GPT-oss-120B critique, DS70B inference)

After reading the scientific publication full text provided above, **extract** the gate, source, drain materials **including any dopant element specified** and the structure-design type that are **explicitly stated** in the text, and generate a JSON file **with fields 'gate', 'source', 'drain', and 'structure_design_type'**.

The 'gate' field must contain the material of the gate electrode (or reference electrode for electrolyte-gated devices), not any functionalizing molecules or catalysts.

gate: the material of the **reference electrode** (or gate electrode for non-electrolyte devices) as described in the paper.

For extended-gate (remote) configurations, the gate material is the base metal of the extended gate electrode itself (e.g., the copper pad on a PCB), not the sensing membrane or functional layer attached to it.

If the paper only mentions a gate-related material such as a catalyst or substrate, use that material for the gate field.

Normalize material names to canonical forms (e.g., "Si substrate with HfO2" -> "phosphorus/silicon"; "Ti/Au" -> "titanium/gold").

If any required field (gate, source, drain, structure_design_type) is not explicitly stated, fill them with defaults: gate = "phosphorus/

```

3780     silicon", source = "gold", drain = "gold", structure_design_type = "
3781     Standard".
3782
3783 **Note:** For electrolyte-gated devices, the gate material is the
3784     reference electrode (e.g., Ag/AgCl; if given by a common name such as
3785     calomel, express it as the chemical pair 'mercury/mercury chloride')
3786     ; for BioFETs the gate material is the functional membrane
3787     composition (e.g., polypyrrole/urease). The source and drain are the
3788     channel contact metals (e.g., Ni/Au).
3789
3790 **Clarification:** The "gate" field refers to the material of the
3791     reference electrode (or gate electrode for solid-state gates). For
3792     back-gated devices, the gate electrode is the substrate; extract its
3793     material from the substrate description (e.g., "n-doped Si" -> "boron
3794     /silicon").
3795
3796 **Include any dopant name preceding the substrate (e.g., "phosphorus/
3797     silicon"); if none is mentioned, use just the substrate material.**
3798
3799 **If the text explicitly states that the substrate serves as the gate (e.
3800     g., "heavily phosphorus-doped silicon wafer (the substrate) serves as
3801     the bottom-gate electrode"), use the substrate material as the gate
3802     .**
3803
3804 **For source and drain, extract the semiconductor channel material of the
3805     FET as explicitly described (e.g., silicon, poly-Si, metal-oxide,
3806     organic semiconductor), not the contact metal unless the
3807     semiconductor itself is not specified.**
3808
3809 If the paper mentions more than one distinct FET configuration (different
3810     gate materials, source/drain metals, or structural designs), create
3811     a separate record for each configuration.
3812
3813 **Structure Design Type** must be selected from the valid values list,
3814     for example:
3815 - *Remote (e.g., the gate electrode is physically separated from the
3816     transistor channel, as in an Extended Gate FET - EGFET)*
3817 - *Electrolyte-Gated (e.g., a gate electrode in direct contact with an
3818     electrolyte)*
3819 - *Top-Gated*
3820 - *Back-Gated*
3821 - *Planar*
3822 - *Vertical*
3823 - *Coplanar*
3824 - *Interdigitated*
3825
3826 **Decision rule:** If the gate electrode is a distinct, physically
3827     separate electrode immersed in the same electrolyte as the channel (i.
3828     e., the gate does not sit directly on the channel surface),
3829     label the 'structure_design_type' as 'Remote'. If the electrolyte
3830     itself serves as the gate dielectric directly contacting the channel
3831     surface (e.g., ion-gel or liquid electrolyte on top of the channel),
3832     label it as 'Electrolyte-Gated'. Use the remaining terms as
3833     defined in the original list.
3834
3835 For instance:
3836 {
3837   "records": [
3838     {
3839       "gate": "silver/silver chloride",
3840       "source": "silver",
3841       "drain": "silver",
3842       "structure_design_type": "Floating"
3843     }
3844   ]
3845 }

```

```

3834         // continue if the publication has recorded multiple different device
3835         architectures
3836     ]
3837 }
3838

```

Material Layer Composition (7 fields) Human initial prompt

3839
3840 After reading the scientific publication full text provided above, please
3841 try to generate formatted JSON file for extraction of important
3842 information. Leave "" for not available fields.

3843
3844 For instance:

```

3845 {
3846   "records": [
3847     {
3848       "substrate": "(substrate material name)",
3849       "substrate_thickness": "xx um",
3850       "channel": "(channel/active layer material name)",
3851       "dielectric_layer": "(dielectric layer material name)",
3852       "dielectric_layer_thickness": "xx nm",
3853       "surface_functionalization": "(surface modification material or
3854       molecule name)",
3855       "structure_dimensionality": "(0D/1D/2D/3D, describing nanostructure
3856       dimensionality)"
3857     }
3858     (continue if the publication has recorded multiple different material
3859     configurations)
3860   ]
3861 }

```

Optimized prompt (M8, GPT-oss-120B critique, DS70B inference)

3861 After reading the scientific publication full text provided above, please
3862 try to generate a formatted JSON file for extraction of important
3863 information.

3864 Leave "" for not available fields, and for `substrate_thickness` and `dielectric_layer_thickness` extract the exact numeric value and unit explicitly given for the substrate or dielectric layer respectively (ignore other layer thicknesses).

3868
3869 **If the substrate thickness is mentioned, always include it in the `substrate_thickness` field, capturing the exact numeric value and unit (e.g., "525 um"). If not stated, apply the following defaults: silicon wafer - 525 um; glass - 1 mm; quartz - 500 um; flexible polymer (e.g., PET) - 125 um; other substrates - leave the field empty.**

3874
3875 **Note: The "channel" field should denote the sensing transistor used in the ISFET system (e.g., p-type Si-FET, organic TFT), not the reference-electrode metal stack.**

3877
3878 When filling the "substrate" field, **select the material that is described as the supporting bulk or wafer on which the device is built; do not treat active sensing materials or coatings as the substrate.**

3881
3882 If the substrate is listed with multiple components separated by "/", retain the full slash-separated string as the substrate name; if the substrate is described as 'degenerately doped silicon' (or any similar phrasing), treat the substrate as plain silicon and do not infer or output the dopant element; otherwise, if the substrate is a doped semiconductor, include the dopant element (e.g., phosphorus-doped silicon) as part of the substrate description.

```

3888 **If the paper mentions multiple substrates, report the substrate that
3889     hosts the transistor channel (the material on which the FET itself is
3890     built) and ignore substrates used only for auxiliary components such
3891     as reference electrodes or packaging.**
3892
3893 **If the paper mentions a generic substrate name such as "glass",
3894     interpret it as a silicon-based substrate and output the material as
3895     "silicon dioxide/sodium oxide/calcium oxide"; assign a default
3896     thickness of "1000 um" unless a specific thickness is explicitly
3897     provided.**
3898
3899 For the `dielectric_layer` field, **only include the material that
3900     functions as the gate insulator in the final device; ignore any
3901     polymer or resist layers that are used solely for processing (e.g.,
3902     PMMA, photo-resist).**
3903
3904 **Only list a dielectric layer if the paper explicitly mentions a
3905     separate gate-dielectric material (e.g., Al2O3, HfO2, Si3N4) distinct
3906     from any native oxide on the substrate.**
3907
3908 If the device contains more than one layer that could be considered a
3909     dielectric, **do not treat the native oxide layer that directly
3910     contacts the channel surface as the dielectric unless it is
3911     explicitly identified as a gate dielectric**; encapsulating polymers
3912     or protective coatings should be placed under *
3913     surface_functionalization* (or omitted if not a functionalization).
3914
3915 **Only fill *surface_functionalization* when a distinct material is
3916     intentionally added to the gate electrode surface that is different
3917     from the electrode/contact material; if the functional element is the
3918     same as the electrode metal or no separate functionalization is
3919     applied, leave the field empty.**
3920
3921 If the substrate description includes an oxide layer (e.g., Si/SiO2) and
3922     its thickness is provided, treat that oxide as the dielectric layer
3923     and populate `dielectric_layer` and `dielectric_layer_thickness`
3924     accordingly.
3925
3926 For instance:
3927 {
3928   "records": [
3929     {
3930       "substrate": "(substrate material name)",
3931       "substrate_thickness": "(xx um if mentioned)",
3932       "channel": "(channel/active layer material name)",
3933       "dielectric_layer": "(dielectric layer material name)",
3934       "dielectric_layer_thickness": "xx nm",
3935       "surface_functionalization": "(surface modification material or
3936         molecule name)",
3937       "structure_dimensionality": "(0D/1D/2D/3D, describing nanostructure
3938         dimensionality)"
3939       // Create one dictionary for each sensing microneedle on the MMNs (
3940         Na+, K+, Ca2+, and pH). The order of the dictionaries does not
3941         matter.
3942     }
3943     (continue if the publication has recorded multiple different material
3944     configurations)
3945   ]
3946 }
3947
3948 /* When a material is given in the paper as an abbreviation (e.g.,
3949     ZrO$2$, PMMA, PMF), output its **full chemical name** in the JSON (e
3950     .g., "zirconium oxide", "poly(methyl methacrylate)", "poly melamine
3951     co-formaldehyde"). Preserve the order of components as they appear (e

```

3942 .g., "zirconium oxide/poly(methyl methacrylate)/poly melamine co-
3943 formaldehyde"). */
3944

3945 **Sensitivity and Response Characteristics (4 fields) Human initial prompt**

3947 After reading the scientific publication full text provided above, please
3948 try to generate formatted JSON file for extraction of important
3949 information. Leave "" for not available fields.

3950 For instance:

```
3951 {
3952 "records": [
3953   {
3954     "response_time": "xx s",
3955     "recovery_time": "xx s",
3956     "sensitivity_numerator": "(the change in signal, e.g. DeltaR,
3957       DeltaI, DeltaG, or absolute values)",
3958     "sensitivity_denominator": "(the reference value, e.g. R0, baseline
3959       current, or gas concentration)"
3959   }
3960   (continue if the publication has recorded multiple different
3961   performance measurements)
3962 ]
3963 }
```

3964 **Optimized prompt (M8, GPT-oss-120B critique, DS70B inference)**

3965 After reading the scientific publication full text provided above, please
3966 generate a formatted JSON file for extraction of important
3967 information.

3968 **If the required values appear inside a table (including ASCII-style
3969 tables), parse the table rows to locate the appropriate numbers for
3970 response time, recovery time, and sensitivity.**

3971 **Extract the response time and recovery time as numeric values in
3972 seconds, rounding to the nearest whole second, and express them as a
3973 number followed by the letter "s" (e.g., "25 s"). If the source
3974 reports these times in minutes (or any other unit), first convert
3975 them to seconds (e.g., 5 min -> 300 s, multiply by 60).**

3976 **When scanning the paper, look for numeric values immediately followed
3977 by terms such as "response time", "recovery time", "rise time", "
3978 settling time", "response", or "recovery" (case-insensitive).** Leave
3979 "" for not-available fields.

3980 **If any required field cannot be found in the provided text, set that
3981 field to an empty string ("") and do not guess or fabricate numbers.
3982 All numeric values must be extracted directly from the supplied
3983 article; do not infer or calculate values that are not explicitly
3984 reported.**

3985 **Only include the fields listed below; do not add title, authors,
3986 journal, or year.** The JSON must contain exactly these fields for
3987 each record: **response_time, recovery_time, sensitivity_numerator,
3988 sensitivity_denominator**.

3989 **If the paper reports an absolute Dirac-voltage shift (e.g., 14 V) and
3990 the corresponding analyte concentration (e.g., 300 pM), compute '
3991 sensitivity_numerator' as the voltage shift converted to millivolts
3992 (1 V = 1000 mV) and set 'sensitivity_denominator' to the reported
3993 concentration using the same unit.**

3994 **Do not enclose the JSON in any markdown or code-fence tags. Output the
3995 raw JSON only.** **Return exactly one JSON object (i.e., a single-

```

3996     element array) summarizing the sensor performance; choose the most
3997     representative values if several are given.**
3998
3999 **The required sensitivity values are located in *Supplementary Table 1*
4000     (the row titled "Our work"). Use the percentage value (e.g., 190 %)
4001     as `sensitivity_numerator` and the corresponding detection-limit
4002     value (e.g., 20 ppm) as `sensitivity_denominator`; ignore any
4003     asterisk footnotes.**
4004
4005 **If multiple sensitivity values appear elsewhere in the text, choose the
4006     one that pertains to the sensor described in the "Current Work"
4007     section (the dEGFET sensor with electropolished Cu electrodes).** **
4008     If several sensitivity percentages are reported for that sensor,
4009     select the highest percentage and use the corresponding concentration
4010     reported with that maximum value.**
4011
4012 **If the paper gives a sensitivity like "<value> mV/pH" (or "<value> mV
4013     per pH unit"), put the number with its unit in `sensitivity_numerator`
4014     and "1 pH" (or "1 dec" for decade) in `sensitivity_denominator`.
4015     Use the primary sensor's value if several are listed.**
4016
4017 **Specifically, when a linear relationship is expressed as a slope (e.g.,
4018     "DeltaV_Dirac = 9.877 mV per decade"), treat the slope value with
4019     its unit as `sensitivity_numerator` and set `sensitivity_denominator`
4020     to "1 dec".**
4021
4022 **For the sensitivity fields in FET gas-sensor papers, locate the
4023     percentage change in drain current reported for the lowest NO2
4024     concentration (e.g., "6.68% for 1 ppm"). Use that percentage as `
4025     sensitivity_numerator` and the corresponding concentration (including
4026     its unit) as `sensitivity_denominator`.**
4027
4028 If the paper reports more than one sensor (different MOF films, gases, or
4029     operating conditions), **output a JSON record for the sensor with
4030     the highest NO2 detection sensitivity, filling in the fields for that
4031     sensor.**
4032
4033 **Pattern rule:** For each required field, find the **first** numeric
4034     value that is **immediately followed** by the appropriate unit--`s`
4035     for `response_time` and `recovery_time`, `%` for `
4036     sensitivity_numerator`, and a concentration unit such as `ppm`, `ppb
4037     `, `ppbv`, `uM`, `mM`, or `M` for `sensitivity_denominator`. Apply
4038     any unit conversions described above. If no such pattern exists, set
4039     the field to an empty string.
4040
4041 **If the denominator is given in log[DA], convert it to ppm** using the
4042     conversion factor provided in the paper; if the paper does not supply
4043     a factor, assume  $1 \log[DA] \approx 0.153 \text{ ppm}$ .
4044
4045 For instance:
4046 {
4047   "records": [
4048     {
4049       "response_time": "xx s",
4050       "recovery_time": "xx s",
4051       "sensitivity_numerator": "(absolute change in drain current DeltaI
4052         per pH unit)",
4053       "sensitivity_denominator": "(the reference value, e.g. R0, baseline
4054         current, or gas concentration). **If expressed as a
4055         concentration (M, mM, uM), convert to ppm before reporting.**"
4056     }
4057   ]
4058 }

```

4050 Synthesis and Thermal Processing (5 fields) Human initial prompt

4051 After reading the scientific publication full text provided above, please
 4052 try to generate formatted JSON file for extraction of important
 4053 information. Leave "" for not available fields.
 4054

4055 For instance:

```
4056 {
4057   "records": [
4058     {
4059       "annealing_temperature": "xx degC",
4060       "annealing_time": "xx h",
4061       "annealing_atmosphere": "(e.g. air, N2, Ar, O2, vacuum, H2)",
4062       "hydrothermal_temperature": "xx degC",
4063       "hydrothermal_time": "xx h"
4064     }
4065     (continue if the publication has recorded multiple different thermal
4066     processing conditions)
4067   ]
4068 }
```

4068 Optimized prompt (M8, GPT-oss-120B critique, DS70B inference)

4069 Here is the JSON file capturing the hydrothermal synthesis details for
 4070 the ZnO nanosheets described in the passage. After reading the
 4071 scientific publication full text provided above, extract any
 4072 explicitly labeled annealing (if present) and hydrothermal heating
 4073 details (temperature, duration, atmosphere) and generate a formatted
 4074 JSON file. If a liquid-phase treatment (e.g., mixed acid, aqueous
 4075 solution) is given with temperature and time, record them as
 4076 hydrothermal_temperature and hydrothermal_time; be sure to extract
 4077 the hydrothermal reaction duration (e.g., "90 degC for 1.5 h") and
 4078 place it in hydrothermal_time. Extract the relevant records and
 4079 format them as JSON, filling each field only when the information is
 present.

4080 For instance:

```
4081 {
4082   "records": [
4083     {
4084       "annealing_temperature": "xx degC",
4085       "annealing_time": "xx h",
4086       "annealing_atmosphere": "(e.g. air, N2, Ar, O2, vacuum, H2)",
4087       "hydrothermal_temperature": "xx degC",
4088       "hydrothermal_time": "xx h"
4089     }
4090     (continue if the publication has recorded multiple different thermal
4091     processing conditions)
4092   ]
4093 }
```

4094 If the manuscript does not mention an annealing step, ****do not infer or
 4095 estimate****; set annealing_temperature, annealing_time, and
 4096 annealing_atmosphere to empty strings (""). Likewise, if no
 4097 hydrothermal step is described, ****do not infer or estimate****; leave
 4098 hydrothermal_temperature and hydrothermal_time empty. If no explicit
 4099 values are found, ****do not infer or estimate****; leave the
 corresponding fields empty (""), and if no annealing or hydrothermal
 step is described, output a list containing a single record where all
 fields are empty strings. ****Return only a valid JSON object (no
 markdown code fences).****

4100
 4101
 4102
 4103