

A TABLE OF CONTENTS

1. Sec. [B](#) quantify the errors in class mapping.
2. Sec. [C](#) show the cost-performance graphs of GPT-4V and GPT-4o,
3. Sec. [D](#) shows the prompt templates used in the experiments.
4. Sec. [E](#) shows the high-level dialogue between the VLM agent and the Verifier.
5. Sec. [F](#) shows the dataset preparation and setup in our experiments.
6. Sec. [G](#) shows additional implementation details in our experiments.
7. Sec. [H](#) shows the additional experimental results. Sec [H.1](#) analyzes the VLM feedback in accuracy and justify that F_1 is a better metric to evaluate feedback quality. Sec. [H.2](#) shows the qualitative results.

B QUANTITATIVE ANALYSIS IN CLASS MAPPING ERRORS

Assessing VLMs’ zero-shot capabilities with close-set vocabularies highlights language ambiguities. In this work, we rely on off-the-shelf sentence embeddings for the class mapping. To quantify errors introduced by mapping model outputs to close-set class labels, we conducted an additional experiment: We sampled 100 raw outputs from LLaVA-1.5 in ADE. A human (one of the authors) evaluated whether the mapping from raw output to class labels, using sentence embeddings, was correct. Table [5](#) Evaluating open-vocabulary models cheaply and automatically remains an open question. Even human evaluators found 10% of the data difficult to map correctly. We have tried to ensure fair comparisons between approaches by maintaining consistent mapping.

Options	Counts
The mapping is correct.	77
The mapping is incorrect and I can provide the correct one.	13
The mapping is incorrect, but it is hard to find a good one from the close set class labels.	10
Total	100

Table 5: **Human studies in quantifying the error in class mapping.**

C PERFORMANCE-COST TRADEOFF

Despite the advances in VLMs’ semantic grounding through self-correction, the identified self-correction trades compute for performance. Fig. [4](#) shows the GPT-4o performance-cost tradeoff in ADE20k.

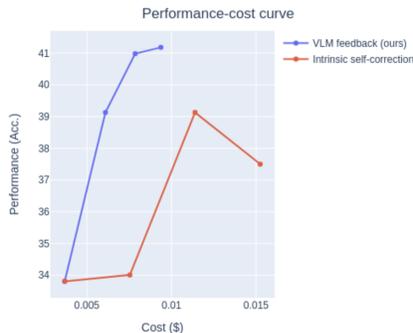


Figure 4: **Cost-performance tradeoff of GPT-4o in ADE20k**

D PROMPT TEMPLATES

We show the full prompt templates

1. To producing base semantic grounding predictions in Fig. 5
2. To enhance previous semantic grounding predictions by taking binary feedback in Fig. 6
3. To enhance previous semantic grounding predictions by taking class label feedback in Fig. 7
4. To produce VLM binary feedback in Fig. 8
5. For the GPT-4V and GPT-4o experiments, we provide the class names by appending ‘*You must answer by selecting from the following names: [COCO or ADE20k Vocabulary]*’ in the prompt¹ as shown in Fig. 9 and Fig. 10.

```
User: You are tasked with visual semantic grounding. Your
→ goal is to determine the class names for objects within a
→ provided image. Each object in the image is identified by
→ a unique ID and its location is defined by a precise
→ bounding box, formatted as: \id{id} \box{[x1, y1, x2,
→ y2]}, where coordinates specify the box corners. The
→ inferred class name for each object is denoted as
→ \class{class name}. Here are the objects: \id{2}
→ \box{[0.1, 0.2, 0.13, 0.43]}
Put your final answer by filling in the placeholder(s) in the
→ following string at the beginning: "\id{2} \box{[0.1,
→ 0.2, 0.13, 0.43]} \class{your answer here}"
```

Figure 5: Prompt template to produce the base predictions. The text in red represents variables.

E EXAMPLE DIALOGUE

In Fig. 11, we demonstrate the iterative interactions between a VLM agent and the Verifier. In Fig. 12, we show the effectiveness of VLM binary verification in GPT-4V².

F DATASET DETAILS

We use ADE20k and COCO panoptic segmentation dataset to evaluate the semantic grounding performance in VLMs. We adopt SoM split provided in the prior work Yang et al. (2023a)³. ADE20k is a large-scale dataset with fine-grained segmentation labels. We adopt the variant with 150 classes, commonly referred to as ADE20k-150. COCO panoptic segmentation is a standard dataset to evaluate visual grounding. There are 133 fine-grained classes in total, composed of 80 thing classes and 53 stuff classes. Consistent with prior works, SoM (Yang et al., 2023a), we use the same subset of 100 images for ADE20k and COCO for our analysis. There are 100 images and 488 segmentation masks in ADE20k SoM split and 101 and 628 segmentation masks in COCO SoM split.

Every region r_i in ADE20k and COCO panoptic segmentation dataset is represented with segmentation mask. We convert them to a more compact representation, *i.e.* bounding box, and feed them to the VLMs in the text prompt

¹<https://github.com/microsoft/SoM/tree/main/benchmark#open-vocab-segmentation-on-coco>

²GPT-4V predictions with simplified prompts as of Mar 22, 2024: <https://imgur.com/a/nbKjIlb>

³<https://github.com/microsoft/SoM/tree/main/benchmark#dataset>

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

```

User: You are tasked with visual semantic grounding. Your
  ↪ goal is to determine the class names for objects within a
  ↪ provided image and leverage the insights from expert
  ↪ analyses. The expert analyses offer detailed information
  ↪ on the inferred class names for each object in the
  ↪ provided image. Each object in the image is identified by
  ↪ a unique ID and its location is defined by a precise
  ↪ bounding box, formatted as: \id{id} \box{[x1, y1, x2,
  ↪ y2]}, where coordinates specify the box corners. The
  ↪ inferred class name for each object is denoted as
  ↪ \class{class name}. I have labeled each object with its
  ↪ ID and overlaid its segmentation mask on the image to
  ↪ clarify the correspondences.

One expert analyses on the provided image are shown below:
* Analysis 1
Object(s) with inferred class names: \id{2} \box{[0.1, 0.2,
  ↪ 0.13, 0.43]} \class{wall}
Expert's decision(s) on class names: The inferred class
  ↪ name(s) for {incorrect obj id} are incorrect. The
  ↪ inferred class name(s) for \id{2} are not "wall".
Expert's suggestion: Adjust the class names for objects with
  ↪ IDs \id{2}

Examine the image and the expert analyses to determine the
  ↪ true class name of the object(s): \id{2} \box{[0.1, 0.2,
  ↪ 0.13, 0.43]}. Put your final answer by filling in the
  ↪ placeholder(s) in the following string at the beginning:
  ↪ "\id{2} \box{[0.1, 0.2, 0.13, 0.43]} \class{your answer
  ↪ here}"

```

Figure 6: Prompt template to improve semantic grounding predictions by taking Binary Feedback. The text in red represents variables.

G IMPLEMENTATION DETAILS

Every experiment throughout this paper is run over three seeds and we report the average scores except for experiments with proprietary VLMs. All the experiments are run in a single-node machine with two A40 GPUs. In the experiments with binary or class label feedback, we only ask VLMs to correct those that are incorrect based on the feedback. Therefore, if the feedback is noisy, *e.g.* VLM binary verification, VLMs can possibly decrease the performances. See Fig. 16 for example.

Open-source VLMs. We adopt LLaVA-1.5 13b (from <https://huggingface.co/llava-hf/llava-1.5-13b-hf>), ViP-LLaVA 13b (from <https://huggingface.co/llava-hf/vip-llava-13b-hf>), and CogVLM (from <https://huggingface.co/THUDM/CogVLM>). When perform the VLM forward pass $o_i = \text{VLM}(x, r_i, q)$, we set the temperature to 0.9, top_p to 0.8, max_new_tokens to 1024, and draw five samples per forward pass. We take the majority vote responses as the final answers o_i .

GPT-4V. As suggested in prior work (Yang et al., 2023a,b), GPT-4V exhibits better grounding ability when the objects are specified by visual prompts rather than text prompts. Therefore, we adopt GPT-4V & SoM to obtain the base predictions, where we overlay object masks and numeric identifiers on the images. Furthermore, when using VLMs to produce feedback, we apply SoM to specify each object. Finally, since GPT-4V has a longer context window compared to open-source VLMs, we include the class list in the prompt to encourage better alignment between the responses and the ground truth. All GPT-4V experiments are done over the OpenAI API and we follow

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

```
User: You are tasked with visual semantic grounding. Your
↪ goal is to determine the class names for objects within a
↪ provided image and leverage the insights from expert
↪ analyses. The expert analyses offer detailed information
↪ on the inferred class names for each object in the
↪ provided image. Each object in the image is identified by
↪ a unique ID and its location is defined by a precise
↪ bounding box, formatted as: \id{id} \box{[x1, y1, x2,
↪ y2]}, where coordinates specify the box corners. The
↪ inferred class name for each object is denoted as
↪ \class{class name}. I have labeled each object with its
↪ ID and overlaid its segmentation mask on the image to
↪ clarify the correspondences.

One expert analyses on the provided image are shown below:
* Analysis 1
Object(s) with inferred class names: \id{2} \box{[0.1, 0.2,
↪ 0.13, 0.43]} \class{wall}
Expert's decision(s) on class names: The inferred class
↪ name(s) for \id{2} are incorrect. The inferred class
↪ name(s) for \id{2} are not "wall".
Expert's suggestion: Adjust the class names for objects with
↪ IDs \id{2} to \class{rail}.

Examine the image and the expert analyses to determine the
↪ true class name of the object(s): \id{2} \box{[0.1, 0.2,
↪ 0.13, 0.43]}. Put your final answer by filling in the
↪ placeholder(s) in the following string at the beginning:
↪ "\id{2} \box{[0.1, 0.2, 0.13, 0.43]} \class{your answer
↪ here}"
```

Figure 7: Prompt template to improve semantic grounding predictions by taking Class Label Feedback. The text in red represents variables.

```
User: Does this cropped image contain "wall"? Answer yes or
↪ no.
```

Figure 8: Prompt template to derive VLM binary feedback. The text in red represents variables.

the exact same evaluation procedures described in Sec. 3.3, where we use the off-the-shelf text embeddings (Huggingface) to map the GPT-4V outputs o_i to the nearest label from the class label list.

We follow the implementation provided in Yang et al. (2023a)⁴ and set the system prompt as: - *For any marks mentioned in your answer, please highlight them with []*. We follow Yang et al. (2023a) to set the alpha parameters in SoM as 0.2 and 0.4 in ADE20k and COCO, respectively. We use the endpoint gpt-4-0125-preview.

GPT-4o. Similar to GPT-4V, we empirically found that SoM prompts improve the base predictions in the semantic grounding tasks in ADE20k. We, therefore, hypothesize that GPT-4o benefits by having SoM prompts. We use the endpoint gpt-4o-2024-05-13.

⁴<https://github.com/microsoft/SoM/blob/main/gpt4v.py>

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

```
User: I have labeled a bright numeric ID at the center for
↳ each visual object in the image. Please enumerate their
↳ names. You must answer by selecting from the following
↳ names: [Class list]
```

Figure 9: Prompt template for GPT-4V and GPT-4o to produce the base predictions. Following prior work (Yang et al., 2023a), we include the full class list in the text prompt. The text in red represents variables.

```
User: You are tasked with visual semantic grounding. Your
↳ goal is to determine the class names for objects within a
↳ provided image and leverage the insights from expert
↳ analyses. The expert analyses offer detailed information
↳ on the inferred class names for each object in the
↳ provided image. Each object in the image is identified by
↳ a unique ID and its location is defined by a precise
↳ bounding box, formatted as: \id{id} \box{[x1, y1, x2,
↳ y2]}, where coordinates specify the box corners. The
↳ inferred class name for each object is denoted as
↳ \class{class name}. I have labeled each object with its
↳ ID and overlaid its segmentation mask on the image to
↳ clarify the correspondences.
```

One expert analyses on the provided image are shown below:

* Analysis 1

Object(s) with inferred class names: \id{2} \box{[0.1, 0.2,

↳ 0.13, 0.43]} \class{wall}

Expert's decision(s) on class names: The inferred class

↳ name(s) for {incorrect obj id} are incorrect. The

↳ inferred class name(s) for \id{2} are not "wall".

Expert's suggestion: Adjust the class names for objects with

↳ IDs \id{2}

Examine the image and the expert analyses to determine the

↳ true class name of the object(s): \id{2} \box{[0.1, 0.2,

↳ 0.13, 0.43]}. Put your final answer by filling in the

↳ placeholder(s) in the following string at the beginning:

↳ "\id{2} \box{[0.1, 0.2, 0.13, 0.43]} \class{your answer

↳ here}"

You must answer by selecting from the following names: [ADE

↳ Class List]

Figure 10: Prompt template for GPT-4V to improve semantic grounding predictions by taking Binary Feedback. Following prior work (Yang et al., 2023a), we include the full class list in the text prompt. The text in red represents variables.

H ADDITIONAL RESULTS

H.1 FEEDBACK ACCURACY DOES NOT STRONGLY CORRELATE WITH SEMANTIC GROUNDING WITH ITERATIVELY SELF-GENERATED FEEDBACK

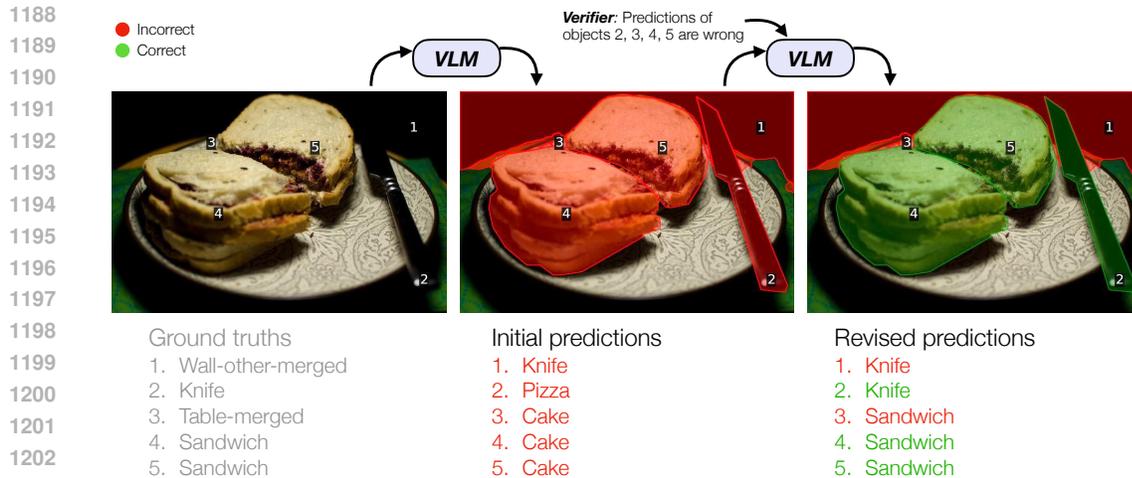


Figure 12: **Enhancing semantic grounding in VLMs with self-generated feedback.** We use GPT-4V as the VLM here. From the left to the center figure, GPT-4V takes the SoM-prompted image [Yang et al. \(2023a\)](#) as input and struggles to predict the class names of each object. From the center to the right figure, GPT-4V takes the same SoM-prompted image and the additional feedback from the verifier and successfully corrects the class names of three out of five objects. The verifier is another GPT-4V that operates on an altered input image and may produce noisy feedback, *e.g.*, misclassify object 1 as correct.

H.2 QUALITATIVE RESULTS

We share additional qualitative results on ADE20k and COCO in [Fig. 13](#), [Fig. 14](#), [Fig. 15](#). We also note that most of the failure cases occur when 1) the VLMs keep their own predictions even though the feedback refers them as incorrect predictions or 2) when the self-generated feedback is incorrect, as shown in [Fig. 16](#).

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

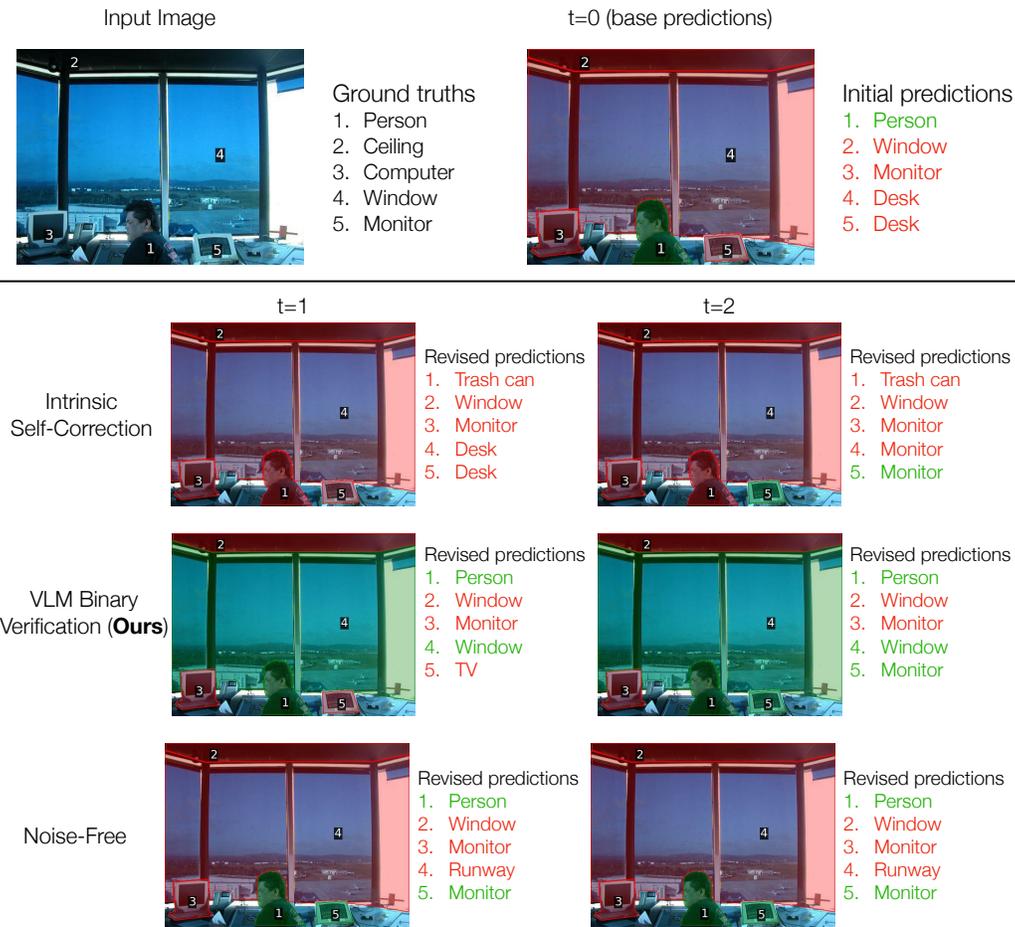


Figure 13: **LLaVA-1.5 qualitative results in ADE20k.** We visualize the predictions of LLaVA-1.5 at time steps from 0 to 2. Intrinsic self-correction fails to identify which predictions are correct/incorrect, while VLM binary verification and Noise-free feedback provide explicit signal on each region, leading to a better chance of correction. From $t = 0$ to $t = 1$, we find that VLM might produce different results (object 4) even when receiving the same feedback (VLM binary verification and Noise-free). As explained in Appendix G in the VLMs forward pass, we draw multiple sequences and take the majority vote as the final responses. For the sake of visualization, we put a bright ID on each object and highlight the incorrect predictions in red and the correct predictions in green.

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

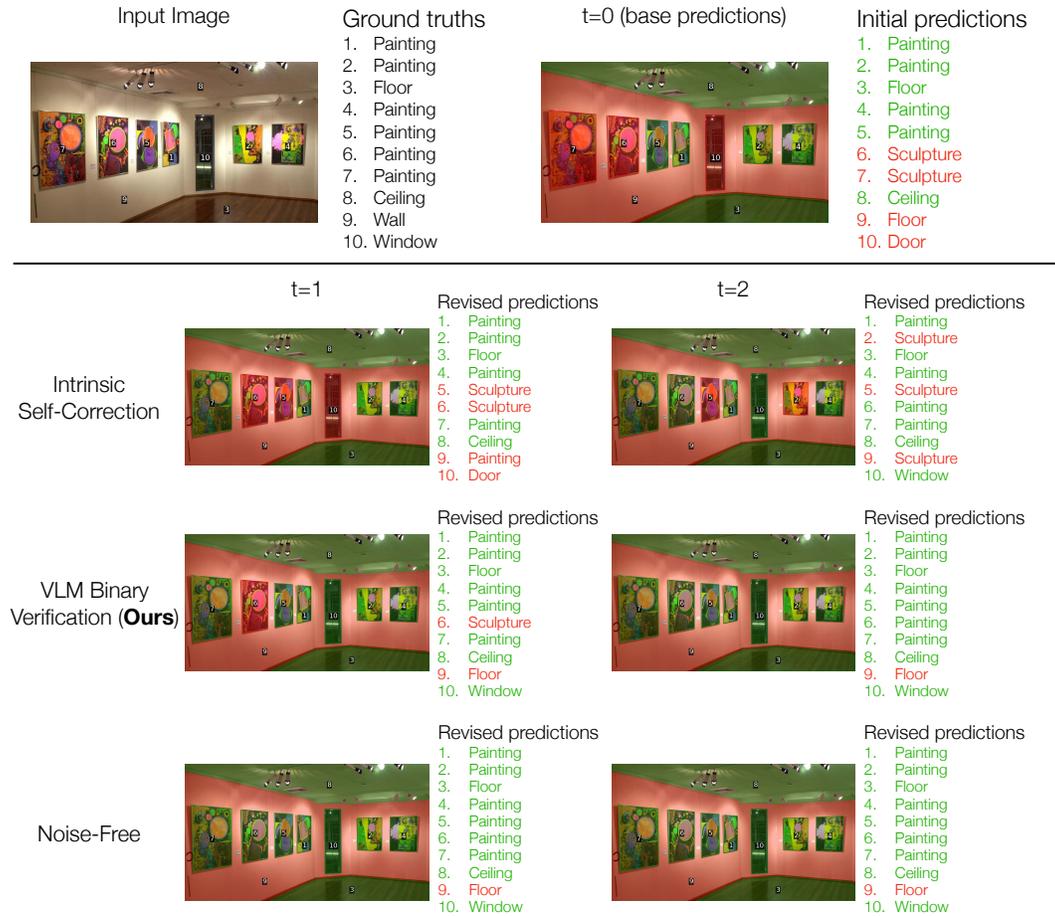


Figure 14: **ViP-LLaVA qualitative results in ADE20k.** We visualize the predictions of ViP-LLaVA at time steps from 0 to 2. Intrinsic self-correction fails to identify which predictions are correct/incorrect, while VLM binary verification and Noise-free feedback provide explicit signal on each region, leading to a better chance of correction. Note that we draw multiple samples in the VLM forward pass, therefore, leading to slightly different results even when the image and query are the same (See Appendix G). For the sake of visualization, we put a bright ID on each object and highlight the incorrect predictions in red and the correct predictions in green.

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

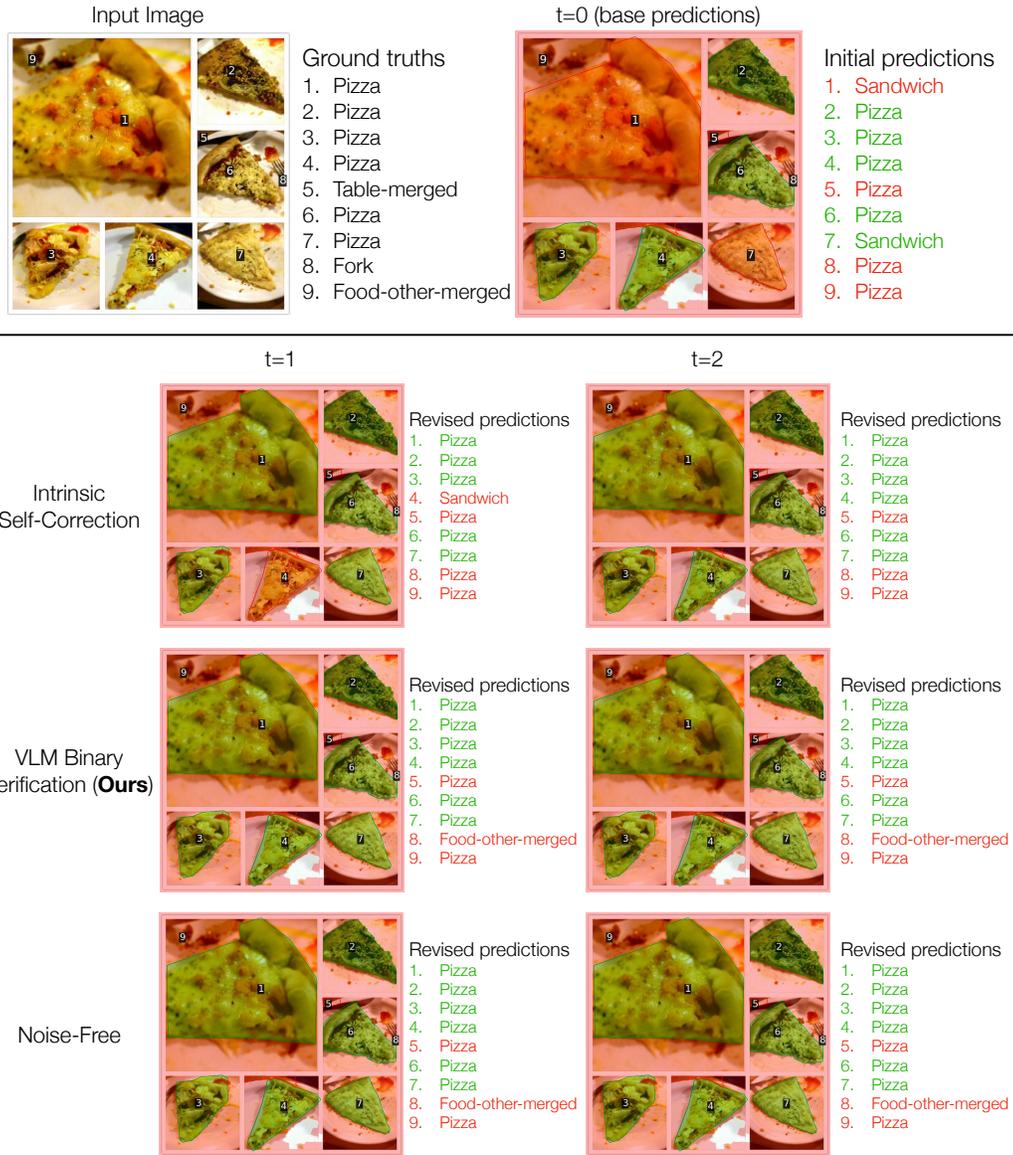


Figure 15: **CogVLM qualitative results in COCO.** We visualize the predictions of CogVLM at time steps from 0 to 2. For the sake of visualization, we put a bright ID on each object and highlight the incorrect predictions in red and the correct predictions in green.

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

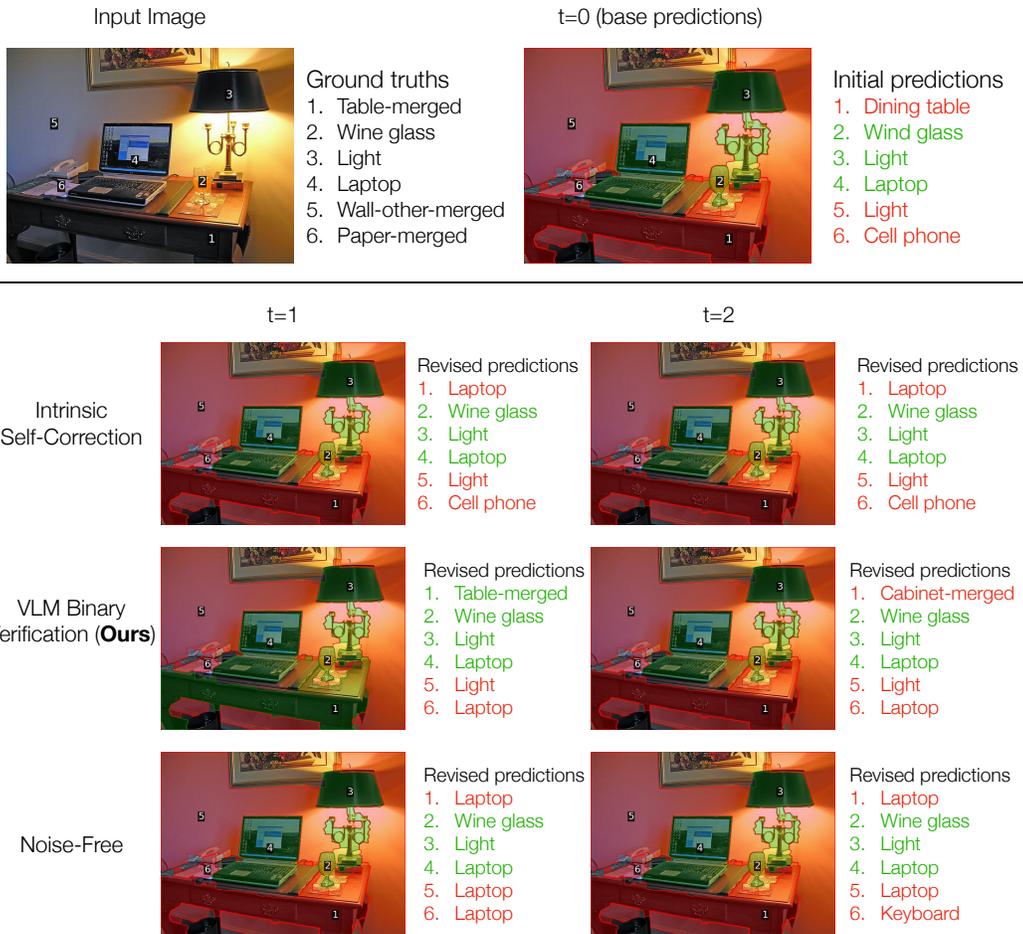


Figure 16: **[Failure case study] LLaVA-1.5 qualitative results in COCO.** All three approaches cannot fix the errors in the initial predictions. For VLM binary verification, from $t = 1$ to $t = 2$, the predictions changes from correct (table-merged) to incorrect (cabinet-merged) since the VLM verifier is not perfect and, therefore, providing misleading feedback. Even with the noise-free feedback, LLaVA-1.5 struggle to adjust the predictions. For the sake of visualization, we put a bright ID on each object and highlight the incorrect predictions in red and the correct predictions in green.