

A. Dataset statistics

Fig. 4 shows the distribution of question types in the Super-CLEVR dataset. The question type is determined by the type of the last operation in the question program.

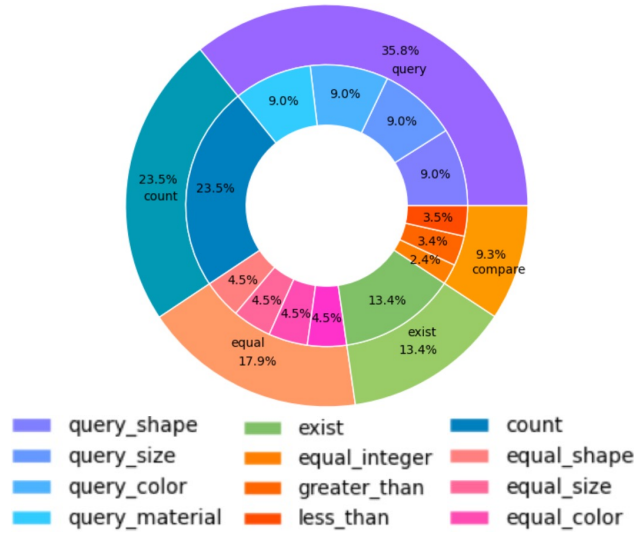


Figure 4. Distribution of question types in Super-CLEVR.

B. List of objects

Super-CLEVR contains 21 objects from 5 categories: airplane, bicycle, bus, car and motorcycle. They are shown in Fig. 5 and Tab. 4.

category	objects
airplane	airliner, biplane, jet, fighter
bicycle	utility bike, tandem bike, road bike, mountain bike
bus	articulated bus, double bus, regular bus, school bus
car	truck, suv, minivan, sedan, wagon
motorcycle	chopper, scooter, cruiser, dirtbike

Table 4. Super-CLEVR dataset contains 21 objects from 5 categories.

C. Dataset controlling

Fig. 6 shows the concept distribution for dataset variants *bal*, *slt* and *long*; and variants *head*, *tail* and *oppo* for testing purpose. Fig. 7 shows the concept co-distribution Matrix M for controlling the concept compositionality (for variants *co-0*, *co-1* and *co-2*). The descriptions for the variants are in Sec. 3.3.

D. Definition of Relative Degrade

Here we describe *Relative Degrade* for each domain shift factors. We use A_j^i to denote the accuracy of model trained using data variant i and tested with data variant j .

Visual complexity:

$$RD = Avg(\frac{A_{easy}^{easy} - A_{mid}^{easy}}{A_{easy}^{easy}}, \frac{A_{easy}^{easy} - A_{hard}^{easy}}{A_{easy}^{easy}}, \frac{A_{mid}^{mid} - A_{easy}^{mid}}{A_{mid}^{mid}}, \frac{A_{mid}^{mid} - A_{hard}^{mid}}{A_{mid}^{mid}}, \frac{A_{hard}^{hard} - A_{easy}^{hard}}{A_{hard}^{hard}}, \frac{A_{hard}^{hard} - A_{mid}^{hard}}{A_{hard}^{hard}})$$

Question redundancy,

$$RD = Avg(\frac{A_{rd-}^{rd-} - A_{rd-}^{rd-}}{A_{rd-}^{rd-}}, \frac{A_{rd-}^{rd-} - A_{rd+}^{rd-}}{A_{rd-}^{rd-}}, \frac{A_{rd-}^{rd-} - A_{rd+}^{rd-}}{A_{rd-}^{rd-}}, \frac{A_{rd-}^{rd-} - A_{rd+}^{rd-}}{A_{rd-}^{rd-}}, \frac{A_{rd+}^{rd-} - A_{rd+}^{rd-}}{A_{rd+}^{rd-}}, \frac{A_{rd+}^{rd-} - A_{rd+}^{rd-}}{A_{rd+}^{rd-}})$$

Concept distribution,

$$RD = \frac{1}{3} \sum_{k \in S} \frac{(A_{head}^k - A_{tail}^k) + (A_{long}^k - A_{oppo}^k)}{2 \cdot A_k^k}, S = \{bal, slt, long\}$$

Concept compositionality,

$$RD = Avg(\frac{A_{co-0}^{co-0} - A_{co-1}^{co-0}}{A_{co-0}^{co-0}}, \frac{A_{co-0}^{co-0} - A_{co-2}^{co-0}}{A_{co-0}^{co-0}}, \frac{A_{co-1}^{co-1} - A_{co-0}^{co-1}}{A_{co-1}^{co-1}}, \frac{A_{co-1}^{co-1} - A_{co-2}^{co-1}}{A_{co-1}^{co-1}}, \frac{A_{co-2}^{co-2} - A_{co-0}^{co-2}}{A_{co-2}^{co-2}}, \frac{A_{co-2}^{co-2} - A_{co-1}^{co-2}}{A_{co-2}^{co-2}})$$

E. More details about P-NSVQA

Given a image containing n objects, we maintain a vector of probability $\mathbf{p} = [p^1, p^2, \dots, p^n]$, where p^k means the probability that object k is selected. We update \mathbf{p} when executing the reasoning operations step by step. In the following, we describe how to compute \mathbf{p} for each kind of operations.

- *scene*
Initialize all the values in \mathbf{p} to 1.
- *filter_identifier[attribute]* (e.g. *filter_color[red]*)
For object k ,

$$p^k = p^k * P_{attribute}^k$$

Here $P_{attribute}^k$ is the probability of object k having the *attribute*, which is predicted by the visual scene parsing model.

- *relate_spacial* (including *relate_behind*, *relate_front*, *relate_right*, *relate_left*)

The output of the *relate* operation is the probabilities of each object being on the *spacial* side of the given object. For example, *relate_left(i)* computes the probabilities of objects to be on the left side of the given object i .

$$p_{front}^k = \frac{1}{1 + e^{-b[(y_k - y_i) + a]}}$$

$$p_{behind}^k = \frac{1}{1 + e^{-b[(y_i - y_k) + a]}}$$

$$p_{right}^k = \frac{1}{1 + e^{-b[(x_k - x_i) + a]}}$$

$$p_{left}^k = \frac{1}{1 + e^{-b[(x_i - x_k) + a]}}$$

Here i is the input object and (x_i, y_i) is the center of it. a and b are hyperparameters. In our experiments, we set $a = 20$ and $b = 0.02$.

- *same_color*, *same_shape*, *same_size*, *same_material*

The *same* operation returns the probabilities of each object having the same attribute as the given object i . For example, for object k and attribute *color*,

$$p^k = cosine_similarity(P_{color}^k, P_{color}^i)$$

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

- *intersect, union*
Given two probability vectors $\mathbf{p}_1, \mathbf{p}_2$, we calculate their intersection or union:

$$\text{Intersection} : \mathbf{p} = \mathbf{p}_1 \odot \mathbf{p}_2$$

$$\text{Union} : \mathbf{p} = 1 - (1 - \mathbf{p}_1) \odot (1 - \mathbf{p}_2)$$

Here \odot is the pointwise product.

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

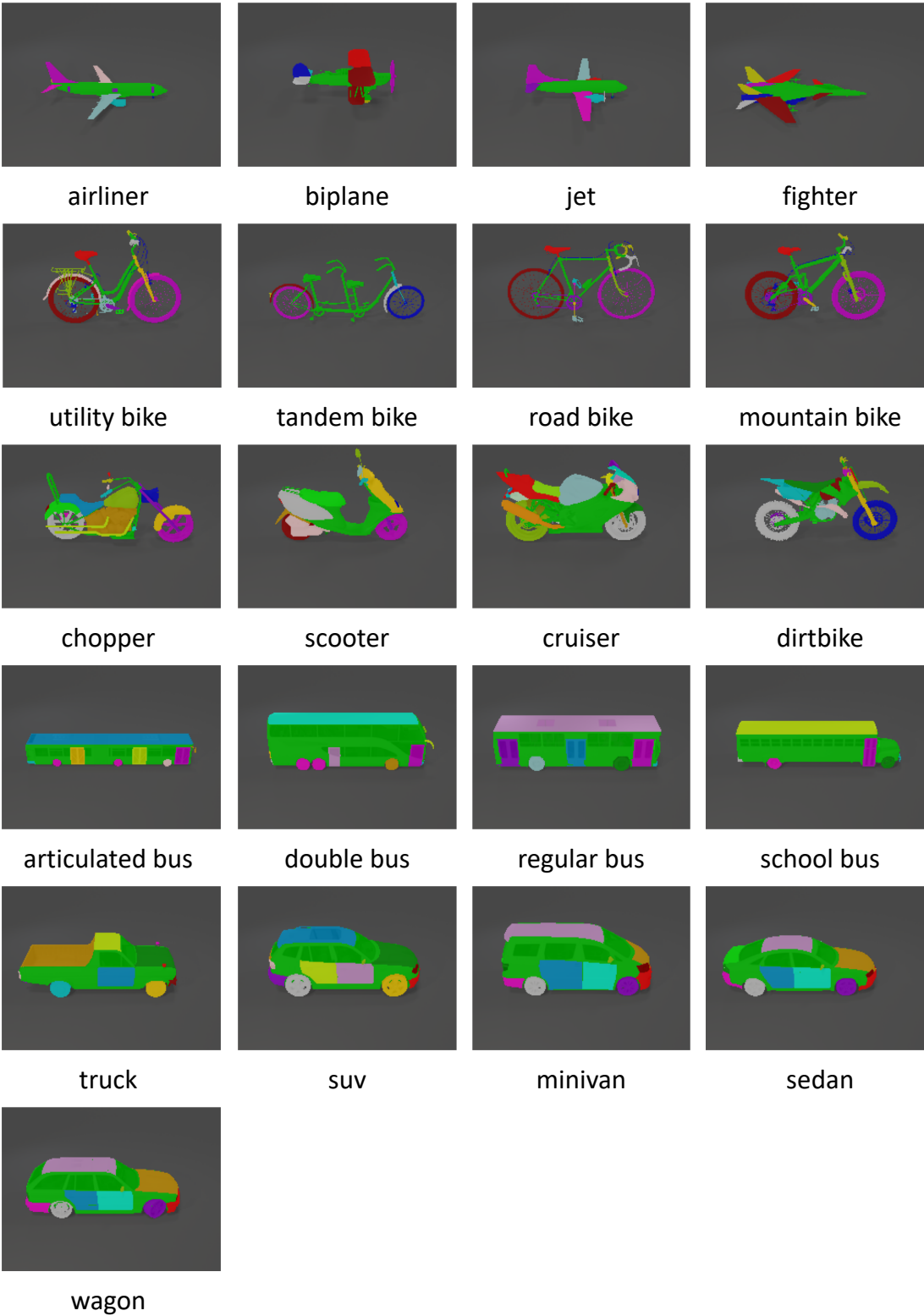


Figure 5. There are 21 objects belonging to 5 categories in the Super-CLEVR dataset.

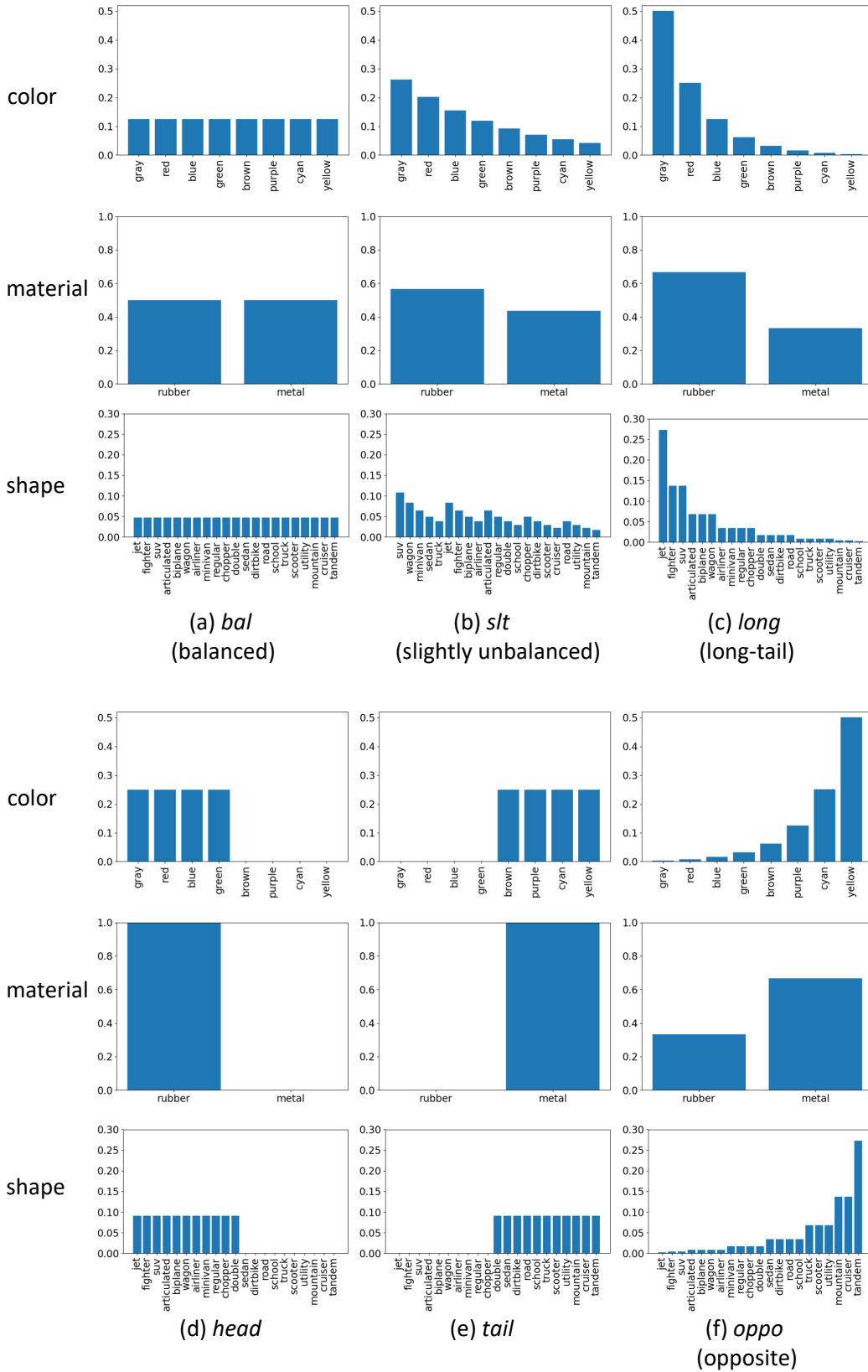


Figure 6. Concept distribution for dataset variants *bal*, *slt*, *long*, *head*, *tail* and *oppo*.

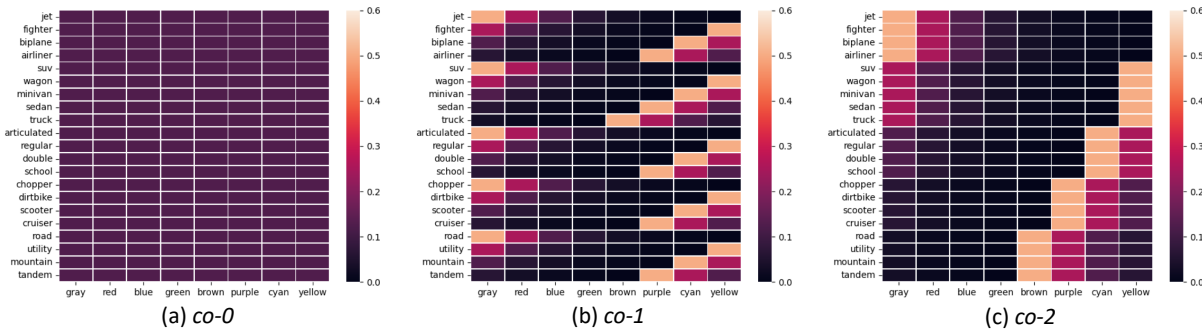


Figure 7. Concept co-distribution matrix M for dataset variants $co-0$, $co-1$ and $co-2$.