

350 **Checklist**

351 1. For all authors...

- 352 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
353 contributions and scope? [Yes]
- 354 (b) Did you describe the limitations of your work? [Yes]
- 355 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Please see
356 Adversarial Example under Section 2
- 357 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
358 them? [Yes]

359 2. If you are including theoretical results...

- 360 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 361 (b) Did you include complete proofs of all theoretical results? [N/A]

362 3. If you ran experiments (e.g. for benchmarks)...

- 363 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
364 mental results (either in the supplemental material or as a URL)? [Yes] The dataset is
365 available in a Google Drive link and the code will be released on a later date.
- 366 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
367 were chosen)? [N/A]
- 368 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
369 ments multiple times)? [Yes]
- 370 (d) Did you include the total amount of compute and the type of resources used (e.g., type
371 of GPUs, internal cluster, or cloud provider)? [N/A]

372 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 373 (a) If your work uses existing assets, did you cite the creators? [Yes] NAO sources images
374 from the OpenImages dataset and is cited.
- 375 (b) Did you mention the license of the assets? [Yes] The annotation is licensed under
376 CC-BY 4.0 and the images are licensed under CC-BY 2.0.
- 377 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 378 (d) Did you discuss whether and how consent was obtained from people whose data you’re
379 using/curating? [N/A]
- 380 (e) Did you discuss whether the data you are using/curating contains personally identifiable
381 information or offensive content? [N/A]

382 5. If you used crowdsourcing or conducted research with human subjects...

- 383 (a) Did you include the full text of instructions given to participants and screenshots, if
384 applicable? [Yes]
- 385 (b) Did you describe any potential participant risks, with links to Institutional Review
386 Board (IRB) approvals, if applicable? [N/A]
- 387 (c) Did you include the estimated hourly wage paid to participants and the total amount
388 spent on participant compensation? [No] The labellers are compensated with money
389 but we are unable to disclose the hourly wage as it is considered trade secrets of Scale
390 AI.

391 A Appendix

392 **Dataset documentation and intended uses** See Appendix B for Datasheet for Dataset of NAO.

393 **URL to website/platform where the dataset/benchmark can be viewed and downloaded by**
394 **the reviewers** The images and dataset can be found on [https://drive.google.com/drive/](https://drive.google.com/drive/folders/15P8s0WoJku6SSEiHLEts860RfytGezi8)
395 [folders/15P8s0WoJku6SSEiHLEts860RfytGezi8](https://drive.google.com/drive/folders/15P8s0WoJku6SSEiHLEts860RfytGezi8)

396 **License** The author bear all responsibility in case of violation of rights and confirmation of the data
397 license.

398 **Hosting, licensing, and maintenance plan** The dataset will be hosted on Google Drive during the
399 review process and will be hosted on scale.com after the review process is completed.

400 **Links to access the dataset and its metadata** The images, annotations and metadata can be found
401 on <https://drive.google.com/drive/folders/15P8s0WoJku6SSEiHLEts860RfytGezi8>

402 **The dataset format** The images for NAO are in either JPEG or PNG format. The bounding box
403 annotation follows the MSCOCO data annotation and can be read using [https://github.com/](https://github.com/cocodataset/cocoapi)
404 [cocodataset/cocoapi](https://github.com/cocodataset/cocoapi).

405 **Long-term preservation** The dataset will be hosted on Google Drive during the review process
406 and will be hosted on scale.com after the review process is completed.

407 **Structured metadata** We will add the metadata (JSON-LD) on the download page of NAO on
408 scale.com once the review process is completed.

409 **Persistent dereferenceable identifier** We will add DOI that directs to the download page of NAO
410 on scale.com once the review process is completed.

411 B Datasheet for Dataset

412 B.1 Motivation

413 **For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific
414 gap that needed to be filled? Please provide a description.

415 NAO is created to evaluate the robustness of object detection models, including out-of-distribution
416 objects. NAO causes many object detection models to incorrectly classify the objects. We hope that
417 NAO can motivate future research on robustness and handling out-of-distribution objects.

418 **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**
419 **company, institution, organization)?**

420 NAO is created by the machine learning team of Scale AI.

421 **Who funded the creation of the dataset?** If there is an associated grant, please provide the name of
422 the grantor and the grant name and number.

423 The ground truth annotation of NAO is funded by Scale AI.

424 B.2 Composition

425 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
426 **countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and
427 interactions between them; nodes and edges)? Please provide a description.

428 Each instance contains the image and the ground truth bounding boxes with their class labels. Some
429 images contain more than one objects.

430 **How many instances are there in total (of each type, if appropriate)?**

Category Name	Count
person	4693
cup	2257
car	752
chair	701
bottle	503
dining table	489
book	315
keyboard	301
bowl	251
fork	243
airplane	214
wine glass	212
dog	210
truck	201
knife	189
cake	178
mouse	176
spoon	163
cell phone	147
potted plant	113
kite	112
handbag	77
tv	76
laptop	75
bicycle	72
sandwich	59
sheep	56
boat	56
bird	55
traffic light	54
backpack	46
couch	41
motorcycle	40
orange	33
sports ball	30
broccoli	29
remote	29
umbrella	29
bench	27
bus	25
vase	25
zebra	24
clock	23

Table 3: Number of objects per object class. Only showing categories with more than 20 objects.

431 There are 7,936 images with 13,604 objects. Please see Table 3 for number of object per category.

432 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of**
433 **instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the sample
434 representative of the larger set (e.g., geographic coverage)? If so, please describe how this representa-
435 tiveness was validated/verified. If it is not representative of the larger set, please describe why not
436 (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

437 NAO does not contain all possible instances. NAO is not a random sample of real world images as it
438 is curated to cause object detection models to produce high confidence errors.

439 **What data does each instance consist of?** “Raw” data (e.g., unprocessed text or images) or features?
440 In either case, please provide a description.

441 Each instance contains an image and ground truth bounding box annotations. The image is in PNG or
442 JPEG format and the bounding box is in MSCOCO (JSON) format to facilitate other researchers to
443 evaluate their models on NAO.

444 **Is there a label or target associated with each instance?** If so, please provide a description.

445 Each image is accompanied by the ground truth bounding boxes.

446 **Is any information missing from individual instances?** If so, please provide a description, ex-
447 plaining why this information is missing (e.g., because it was unavailable). This does not include
448 intentionally removed information, but might include, e.g., redacted text.

449 There is no information missing from individual instances.

450 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social
451 network links)?** If so, please describe how these relationships are made explicit.

452 The images might be taken by the same person.

453 **Are there recommended data splits (e.g., training, development/validation, testing)?** If so,
454 please provide a description of these splits, explaining the rationale behind them.

455 It is recommended that NAO to be used exclusively as a test set. However, a train, validation and test
456 set split is available for methods that require training. There are 7,196 images in the train set, 373 in
457 the validation set, and 367 in the test set.

458 **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a
459 description.

460 Annotation errors are minimized through reviewers and quizzes that are randomly inserted between
461 real annotation tasks. Labellers with quiz score less than a certain threshold were disabled from
462 working on NAO.

463 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,
464 websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there guarantees
465 that they will exist, and remain constant, over time; b) are there official archival versions of the
466 complete dataset (i.e., including the external resources as they existed at the time the dataset was
467 created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources
468 that might apply to a future user? Please provide descriptions of all external resources and any
469 restrictions associated with them, as well as links or other access points, as appropriate.

470 NAO sources images from OpenImages. We provide metadata to link each image in NAO to the
471 original image from OpenImages.

472 **Does the dataset contain data that might be considered confidential (e.g., data that is protected
473 by legal privilege or by doctor-patient confidentiality, data that includes the content of individ-
474 uals’ non-public communications)?** If so, please provide a description.

475 No. All images are derived from OpenImages.

476 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,
477 or might otherwise cause anxiety?** If so, please describe why.

478 No. All images are derived from OpenImages.

479 **Does the dataset relate to people?** If not, you may skip the remaining questions in this section.

480 While NAO does not heavily focus on people, "Person" is one of the main object categories in a
481 universal dataset (OpenImages).

482 **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how
483 these subpopulations are identified and provide a description of their respective distributions within
484 the dataset.

485 No subpopulation of people is identified.

486 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indi-**
487 **rectly (i.e., in combination with other data) from the dataset?** If so, please describe how.

488 The people in the frame might be identified by their Flickr account using the metadata from OpenIm-
489 ages.

490 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that re-**
491 **veals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union**
492 **memberships, or locations; financial or health data; biometric or genetic data; forms of gov-**
493 **ernment identification, such as social security numbers; criminal history)?** If so, please provide
494 a description.

495 The dataset does not contain any sensitive data.

496 **B.3 Collection Process**

497 **How was the data associated with each instance acquired?** Was the data directly observable (e.g.,
498 raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived
499 from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If data was
500 reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If
501 so, please describe how.

502 The images and bounding boxes are directly observable

503 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or**
504 **sensor, manual human curation, software program, software API)?** How were these mechanisms
505 or procedures validated?

506 Please see Section 3 in the main paper for a description of the data collection process.

507 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
508 **probabilistic with specific sampling probabilities)?**

509 The dataset is sampled from OpenImages and images with egregious errors from the EfficientDet-d7
510 are included in NAO.

511 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
512 **and how were they compensated (e.g., how much were crowdworkers paid)?**

513 The creation of NAO include the authors and crowdworkers. The author wrote the instructions,
514 designed the annotation tasks and created the quizzes. The crowdworkers identified if the model pre-
515 dictions are indeed incorrect and put bounding boxes around objects. Crowdworkers are compensated
516 with money but we are unable to disclose the amount as it is considered trade secret.

517 **Over what timeframe was the data collected?** Does this timeframe match the creation timeframe
518 of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please
519 describe the timeframe in which the data associated with the instances was created.

520 It is unknown what timeframe the images were collected but ground truth annotations were collected
521 in October 2020 and May 2021.

522 **Were any ethical review processes conducted (e.g., by an institutional review board)?** If so,
523 please provide a description of these review processes, including the outcomes, as well as a link or
524 other access point to any supporting documentation.

525 We ensured our work conforms to Neurips ethics review guidelines – potential negative societal
526 impacts (Section 2) and general ethical conduct (Section 3).

527 **Does the dataset relate to people?** If not, you may skip the remainder of the questions in this
528 section.

529 The dataset is not directly related to people but people might be included in the images.

530 **Did you collect the data from the individuals in question directly, or obtain it via third parties**
531 **or other sources (e.g., websites)?**

532 The images were sourced from OpenImages which in turns come from Flickr.

533 **Were the individuals in question notified about the data collection?** If so, please describe (or
534 show with screenshots or other information) how notice was provided, and provide a link or other
535 access point to, or otherwise reproduce, the exact language of the notification itself.

536 It is unknown whether the creators of OpenImages informed the image uploaders.

537 **Did the individuals in question consent to the collection and use of their data?** If so, please
538 describe (or show with screenshots or other information) how consent was requested and provided,
539 and provide a link or other access point to, or otherwise reproduce, the exact language to which the
540 individuals consented.

541 Unknown.

542 **If consent was obtained, were the consenting individuals provided with a mechanism to revoke
543 their consent in the future or for certain uses?** If so, please provide a description, as well as a link
544 or other access point to the mechanism (if appropriate).

545 Not applicable.

546 **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data
547 protection impact analysis) been conducted?** If so, please provide a description of this analysis,
548 including the outcomes, as well as a link or other access point to any supporting documentation.

549 Not applicable.

550 **B.4 Preprocessing/cleaning/labeling**

551 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,
552 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing
553 of missing values)?** If so, please provide a description. If not, you may skip the remainder of the
554 questions in this section.

555 No extra processing is done to the images. For the bounding box labeling process, please see Section
556 3 in the main paper.

557 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support
558 unanticipated future uses)?** If so, please provide a link or other access point to the “raw” data.

559 Not applicable.

560 **Is the software used to preprocess/clean/label the instances available?** If so, please provide a
561 link or other access point.

562 No.

563 **B.5 Uses**

564 **Has the dataset been used for any tasks already?** If so, please provide a description.

565 No.

566 **Is there a repository that links to any or all papers or systems that use the dataset?** If so, please
567 provide a link or other access point.

568 Not applicable.

569 **What (other) tasks could the dataset be used for?**

570 NAO can be used to evaluate the robustness of object detection model against out-of-distribution
571 objects.

572 **Is there anything about the composition of the dataset or the way it was collected and prepro-
573 cessed/cleaned/labeled that might impact future uses?** For example, is there anything that a future
574 user might need to know to avoid uses that could result in unfair treatment of individuals or groups
575 (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g., financial harms, legal
576 risks) If so, please provide a description. Is there anything a future user could do to mitigate these
577 undesirable harms?

578 There is minimal risk because the images are already available in public from OpenImages.
579 **Are there tasks for which the dataset should not be used?** If so, please provide a description.
580 No.

581 **B.6 Distribution**

582 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** If so, please provide a description.

583 The dataset is distributed through Google Drive and will be distributed on scale.com

584 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does the dataset have a digital object identifier (DOI)?

585 NAO is initially distributed through Google Drive and later distributed on scale.com once the paper review process is completed. The images are downloadable as a zip file on Google Drive and will be available as a tarball archive on scale.com. The bounding box annotations are available in MSCOCO format. DOI will be provided once the dataset is hosted on scale.com

586 **When will the dataset be distributed?**

587 The dataset is available now on Google Drive and will be distributed on scale.com once the paper review process is completed.

588 **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.

589 The images copyright belongs to the original Flickr users. The bounding box annotations are made available under the CC-BY 4.0 license.

590 **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.

591 Not applicable.

592 **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.

593 Not applicable.

594 **B.7 Maintenance**

595 **Who is supporting/hosting/maintaining the dataset?**

596 NAO is hosted on Google Drive and will be hosted on scale.com once the paper review process is completed. The maintenance will be done by the authors.

597 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

598 All questions and comments can be sent to Felix Lau: felix.lau@scale.com

599 **Is there an erratum?** If so, please provide a link or other access point.

600 All changes to the dataset will be announced in the CHANGELOG file.

601 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** If so, please describe how often, by whom, and how updates will be communicated to users (e.g., mailing list, GitHub)?

602 The dataset might be updated to correct labeling errors and will be communicated using the CHANGELOG file.

622 **If the dataset relates to people, are there applicable limits on the retention of the data associ-**
623 **ated with the instances (e.g., were individuals in question told that their data would be retained**
624 **for a fixed period of time and then deleted)?** If so, please describe these limits and explain how
625 they will be enforced.

626 Not applicable.

627 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please
628 describe how. If not, please describe how its obsolescence will be communicated to users.

629 We plan to only support and host the latest version of NAO.

630 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
631 **them to do so?** If so, please provide a description. Will these contributions be validated/verified? If
632 so, please describe how. If not, why not? Is there a process for communicating/distributing these
633 contributions to other users? If so, please provide a description.

634 Yes. Please email all suggestions or contributions to Felix Lau: felix.lau@scale.com