



Exceptional service in the national interest

But it Looks so Real! Challenges in Training Models with Synthetic Data for International Safeguards

Zoe N. Gastelum, Timothy M. Shead, Matthew
R. Marshall

INSTITUTE OF NUCLEAR MATERIALS MANAGEMENT

ANNUAL MEETING

July , 2022



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc. for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.



Four Challenges:

- 1) Image backgrounds have an unexpectedly large impact on model performance.
- 2) Negative examples are more effective when they include distractors.
- 3) Object configuration and positioning influence identification.
- 4) Computer vision models are learning the wrong features from training data.

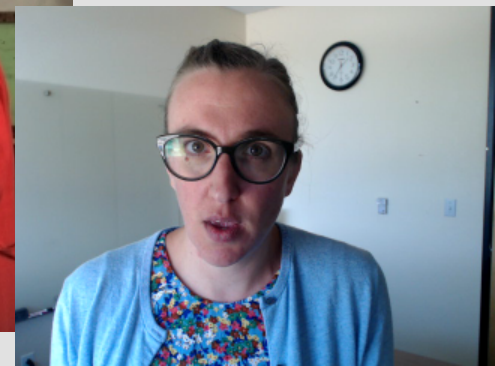
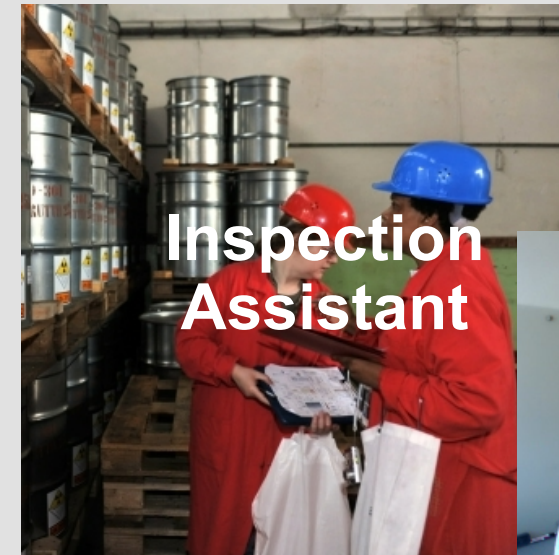
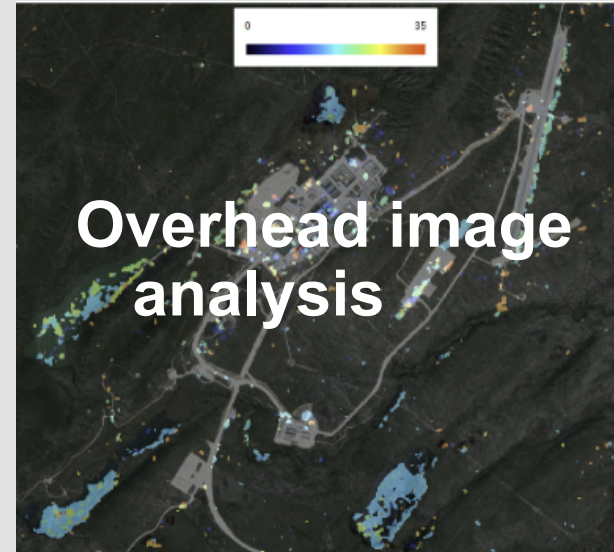




Brief Introduction



Deep learning is being explored for multiple safeguards applications, especially for visual tasks.





We are developing an open dataset for safeguards computer vision research and development.

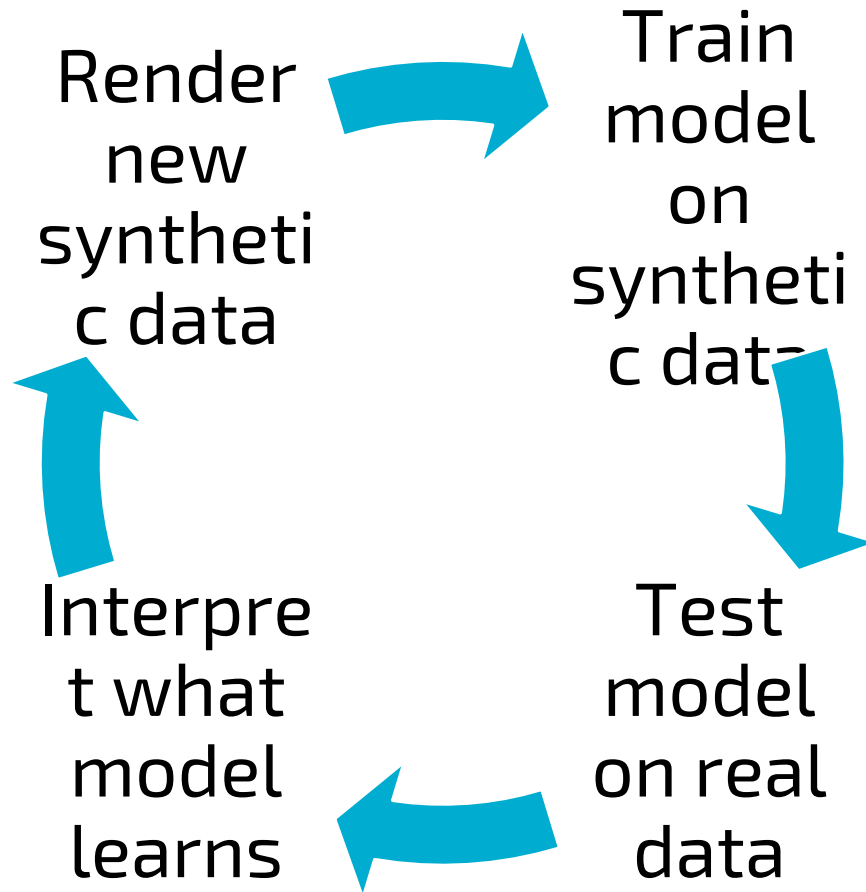
- 2-dimension images of computer models of 30B, 48X, 48Y, and 48G containers
- Vary container size, orientation, lighting, configuration
- Real-world panoramic HDR or synthetic background
- 720K synthetic currently available on Berkeley Data Cloud
- Check out: <https://limbo-ml.readthedocs.io/>



A small, real-world image set was curated for validation.



Our data validation process supports iterative image validation and rendering.



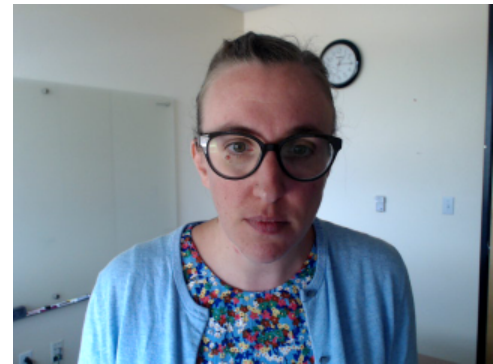
Example Impacts:

1. Relevant containers in rows
2. Partially obscure relevant containers
3. Render synthetic distractors
4. Render synthetic distractors in groups

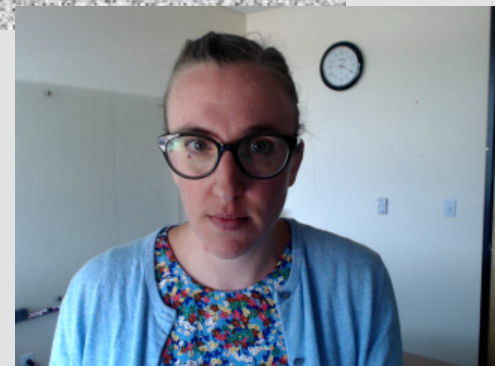
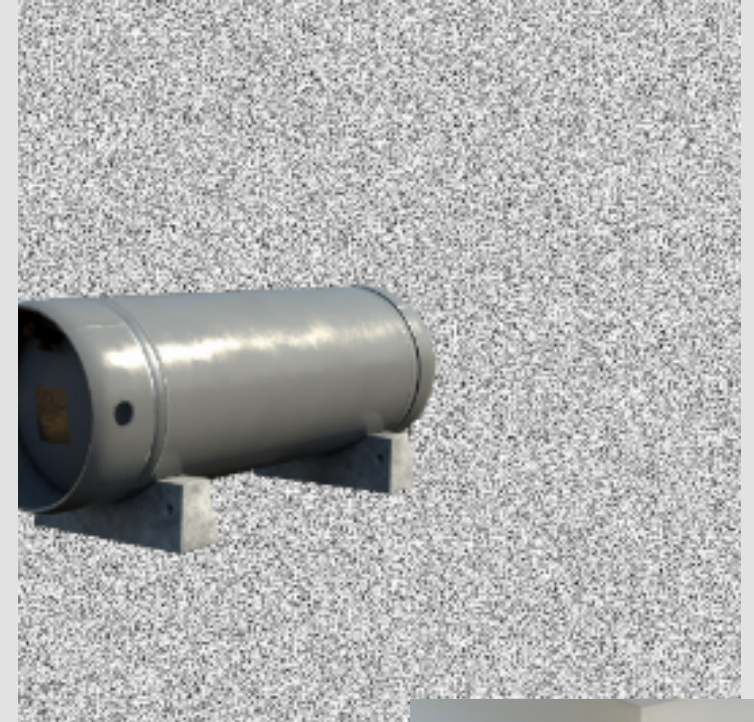
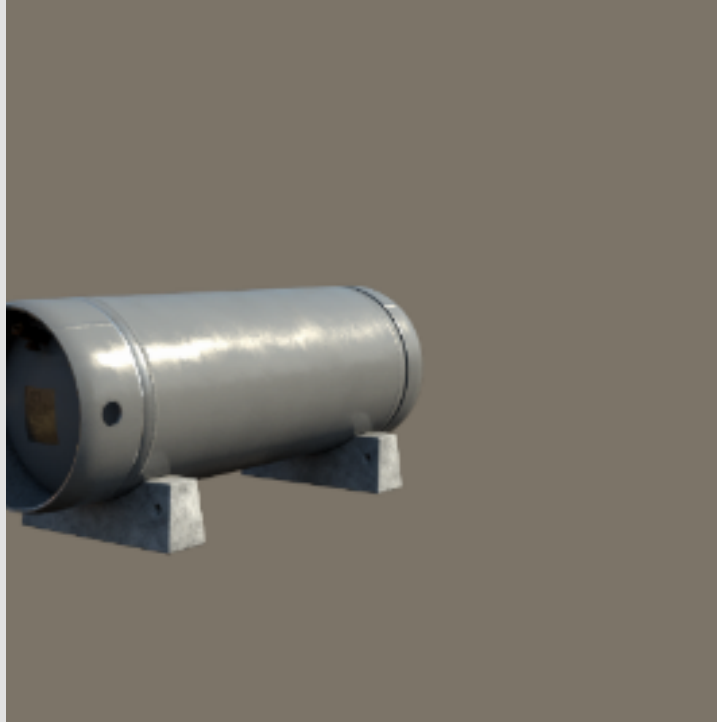
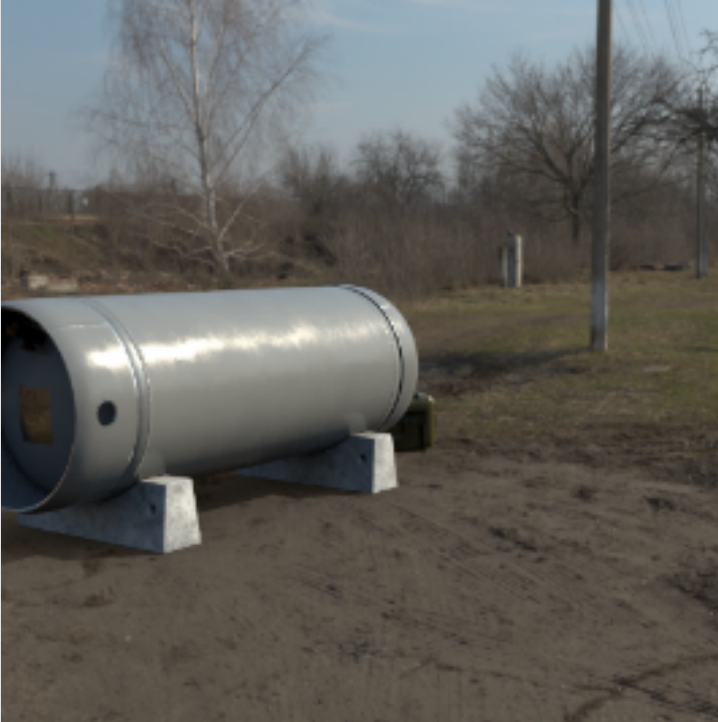




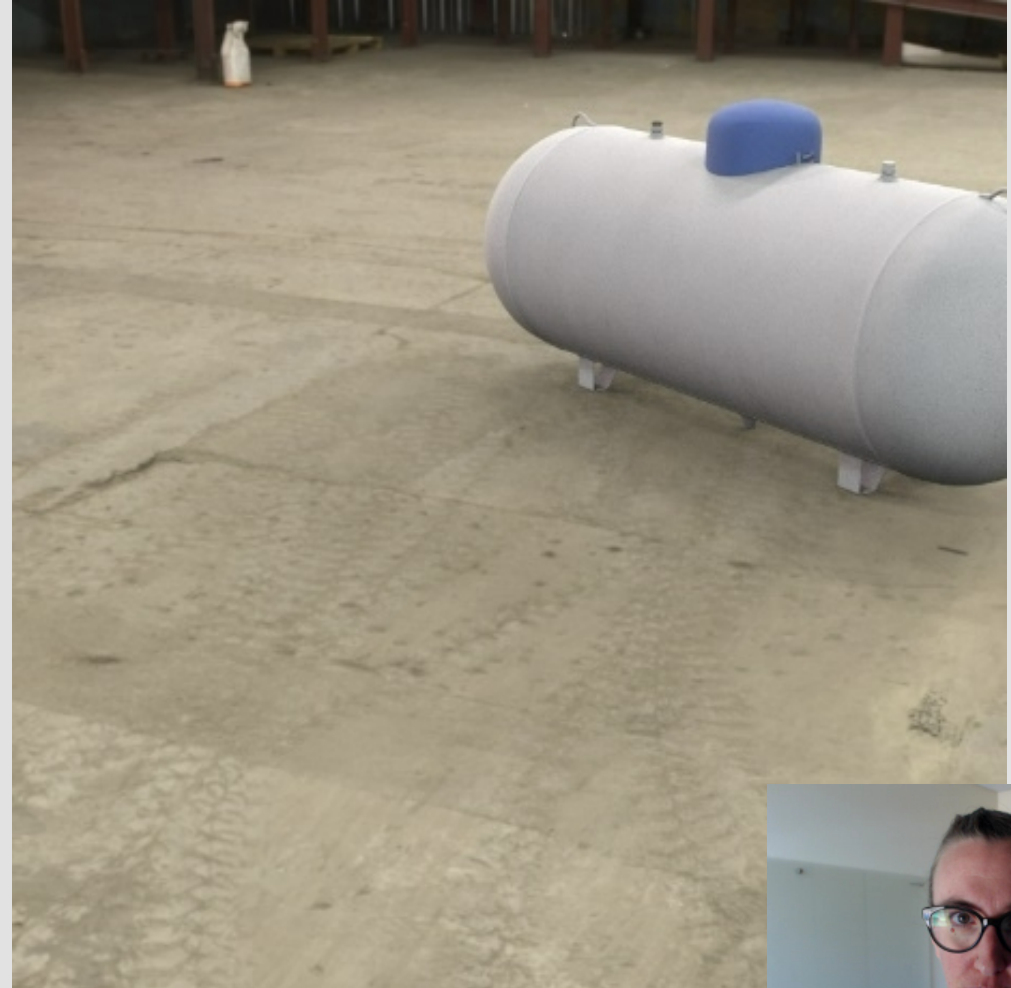
Four Observations



Observation 1: Image backgrounds have an unexpectedly large impact on model performance.



Observation 2: Negative examples are more effective when they include distractors.



Observation 3: Object configuration and positioning influence identification.



Observation 4: Computer vision models are learning the wrong features from training data.





Research Priorities



How might we improve computer vision model training to ensure the relevant data is being learned?

- Irrelevant differences in synthetic images are straining the performance of computer vision models when making inferences on real-world data.
- Adding more training data is not an option.
- International safeguards requires robust inferences.
- We have modified and updated our synthetic images. It is time to re-examine computer vision model training.



Thank you for your attention!

Contact: zgastel@sandia.gov

Data: <https://limbo-ml.readthedocs.io/>
<https://bdc.lbl.gov/>

The National Nuclear Security Administration's Defense Nuclear Nonproliferation R&D sponsored this research. Thank you!

