

Research Review 2021

Knowing When You Don't Know: AI Engineering in an Uncertain World

November 2021

Eric Heim, AI Division

Copyright 2021 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702 -15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

Carnegie Mellon® is registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM21-0875

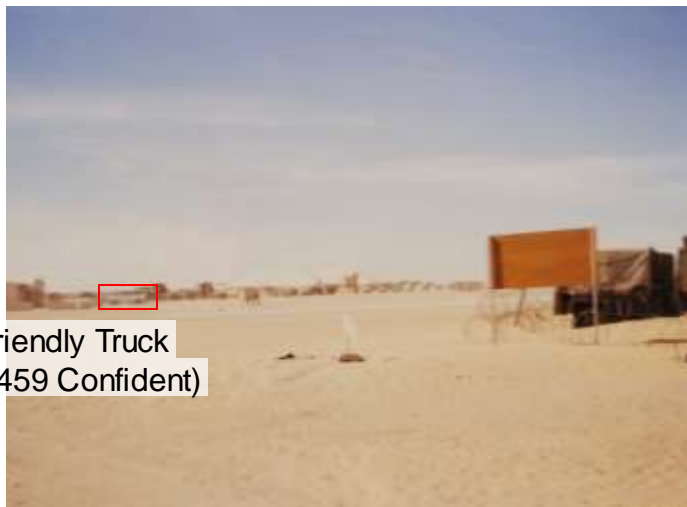
Quantifying Uncertainty: A Key Component for **informative** and Robust AI Systems



Friendly Truck
(0.9834 Confident)

Image: South Carolina National Guard, 151st Signal Battalion

Quantifying Uncertainty: A Key Component for **informative** and Robust AI Systems

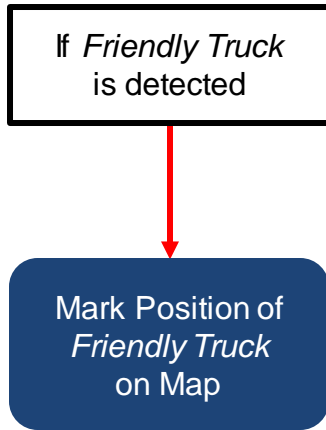


Friendly Truck
(0.3459 Confident)

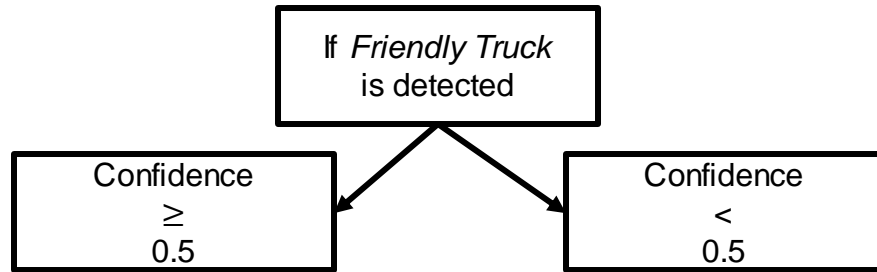
Image: South Carolina National Guard, 151st Signal Battalion

Accurate estimates of uncertainty can lead to better informed **decision making.**

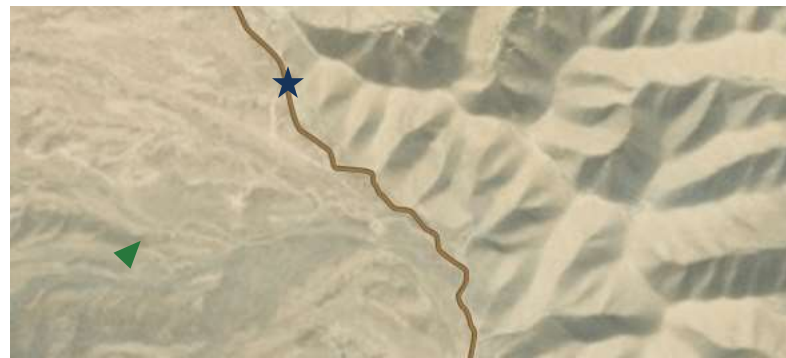
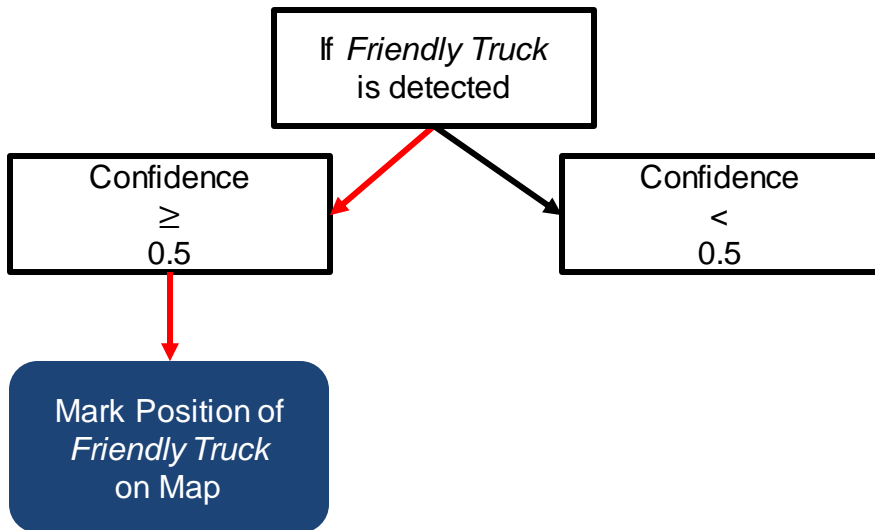
Quantifying Uncertainty: A Key Component for Informative and **Robust** AI Systems



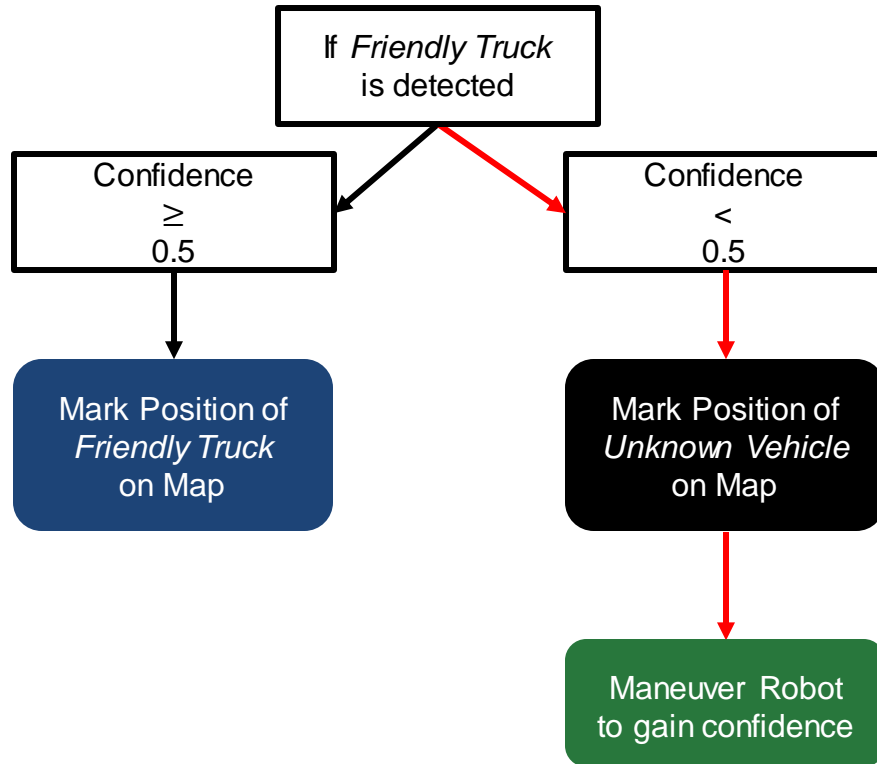
Quantifying Uncertainty: A Key Component for Informative and **Robust** AI Systems



Quantifying Uncertainty: A Key Component for Informative and **Robust** AI Systems



Quantifying Uncertainty: A Key Component for Informative and **Robust** AI Systems



By allowing high-level reasoning to be informed by predictive uncertainty, AI systems can be **more robust** to failures caused by unconfident predictions.

Quantifying Uncertainty: A Key Component for Informative and **Robust** AI Systems

Our Work: Evaluating, Characterizing, Articulating, and Rectifying Uncertainty in ML models for the purpose of more informative and robust AI Systems

Confidence ≥ 0.5

Mark Position
Friendly Truck
on Map

This Talk: Evaluating ML model *calibration*.

Maneuver Robot
to gain confidence

g to be
nty, AI
failures
ions.

Calibration: A Way to Interpret Model Uncertainty



Friendly Truck



Enemy Tank

Classifier

Calibration: A Way to Interpret Model Uncertainty



Friendly Truck



Enemy Tank



Calibration: A Way to Interpret Model Uncertainty



Friendly Truck



Enemy Tank



Calibration: A Way to Interpret Model Uncertainty



Friendly Truck



Enemy Tank



Classifier Calibration: Classifier outputs match the frequency of class labels.

Calibration: A Way to Interpret Model Uncertainty



Friendly Truck



Enemy Tank



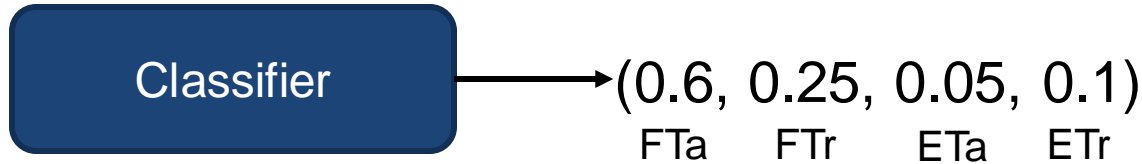
For all possible inputs that the classifier outputs $(0.6, 0.4)$...
 60% of the inputs should be a friendly truck,
 40% of the inputs should be an enemy tank.

Classifier Calibration: Classifier outputs match the frequency of class labels.

Evaluating Classifier Calibration

Modern machine learning literature has focused on evaluating classifier calibration according to their **Top-1 Expected Calibration Error (ECE)**

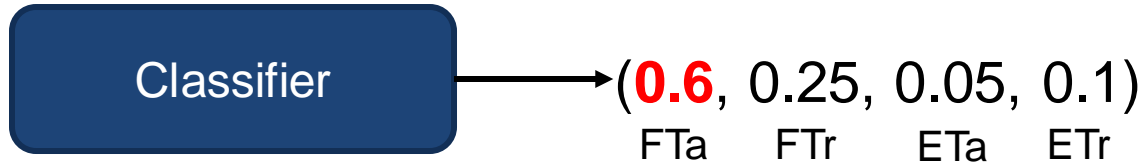
Classes = {Friendly Tank, Friendly Truck, Enemy Tank, Enemy Truck}



Evaluating Classifier Calibration

Modern machine learning literature has focused on evaluating classifier calibration according to their **Top-1 Expected Calibration Error (ECE)**

Classes = {Friendly Tank, Friendly Truck, Enemy Tank, Enemy Truck}



Top-1 Expected Calibration Error (ECE)

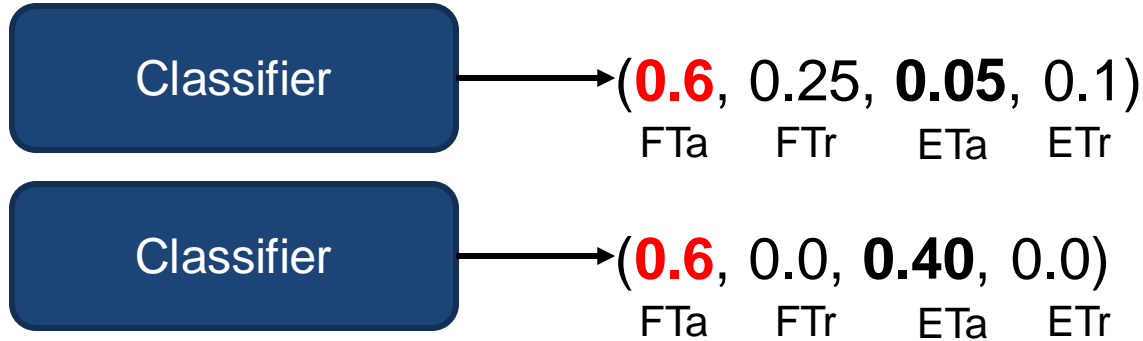
Considers only the most confident class in evaluating for calibration

For **all possible inputs** that the **classifier** outputs **0.6 as the most confident class...**
60% of those inputs should be that class.

Evaluating Classifier Calibration

Modern machine learning literature has focused on evaluating classifier calibration according to their **Top-1 Expected Calibration Error (ECE)**.

Classes = {Friendly Tank, Friendly Truck, Enemy Tank, Enemy Truck}



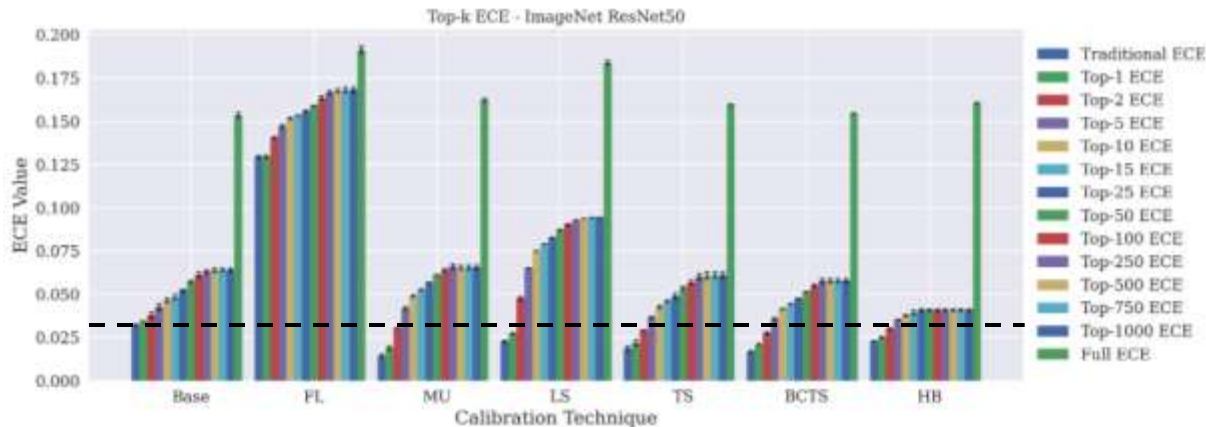
According to Top-1 ECE, these two classifiers ***are considered the same***. However, the two outputs can mean very different things with ***mission context***.

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes



Traditional ECE – Measures miscalibration for most confident class

Overall Goal: Evaluate the state of the art in classifier calibration according to context focused metrics to observe how they perform in different definitions of reliability.

Experiment #1: Top- k

- Data Set: ImageNet
- Base Model: ResNet50

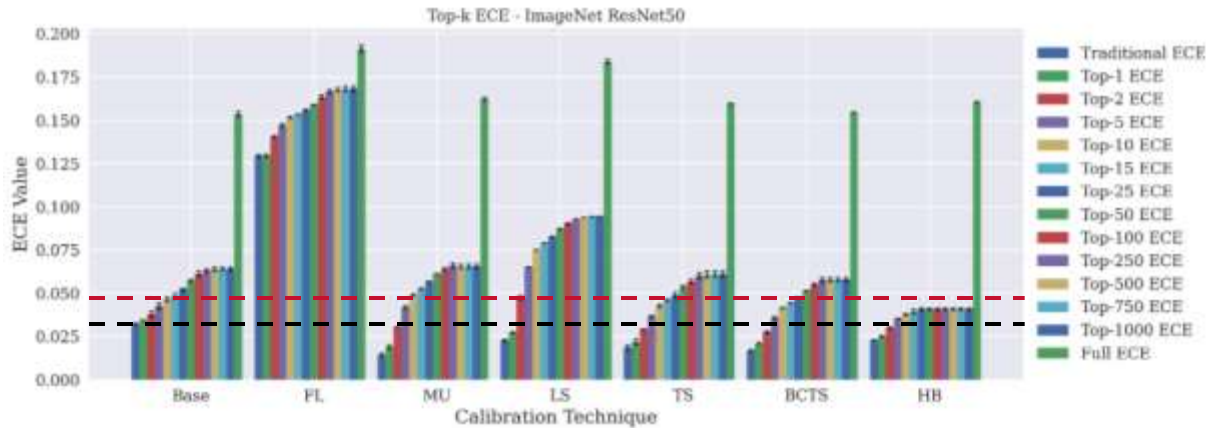
Question: How to these methods perform outside of the most confident class?

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes



Top-10 ECE– Measures miscalibration for the 10 most confident classes

Overall Goal: Evaluate the state of the art in classifier calibration according to context focused metrics to observe how they perform in different definitions of reliability.

Experiment #1: Top- k

- Data Set: ImageNet
- Base Model: ResNet50

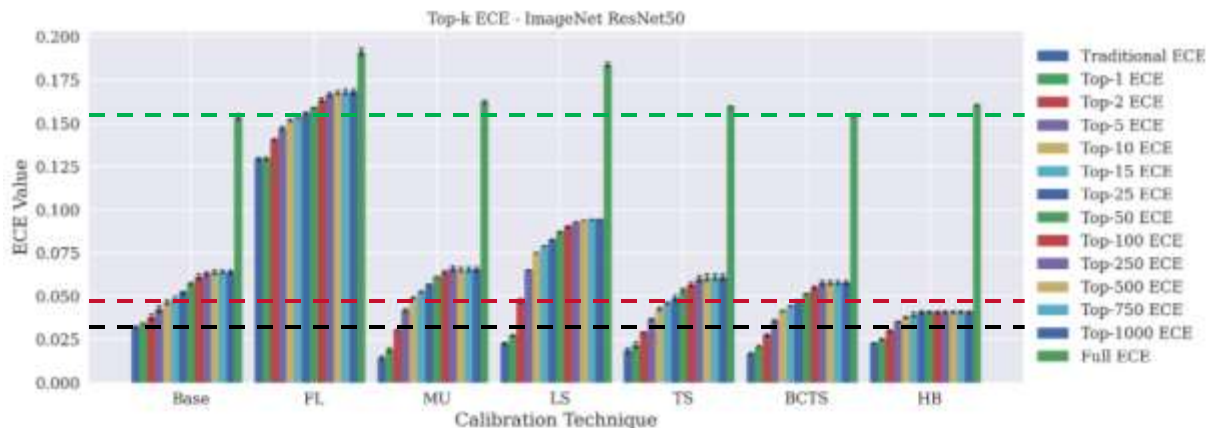
Question: How to these methods perform outside of the most confident class?

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes



Full ECE—Measures miscalibration across all classes

Overall Goal: Evaluate the state of the art in classifier calibration according to context focused metrics to observe how they perform in different definitions of reliability.

Experiment #1: Top- k

- Data Set: ImageNet
- Base Model: ResNet50

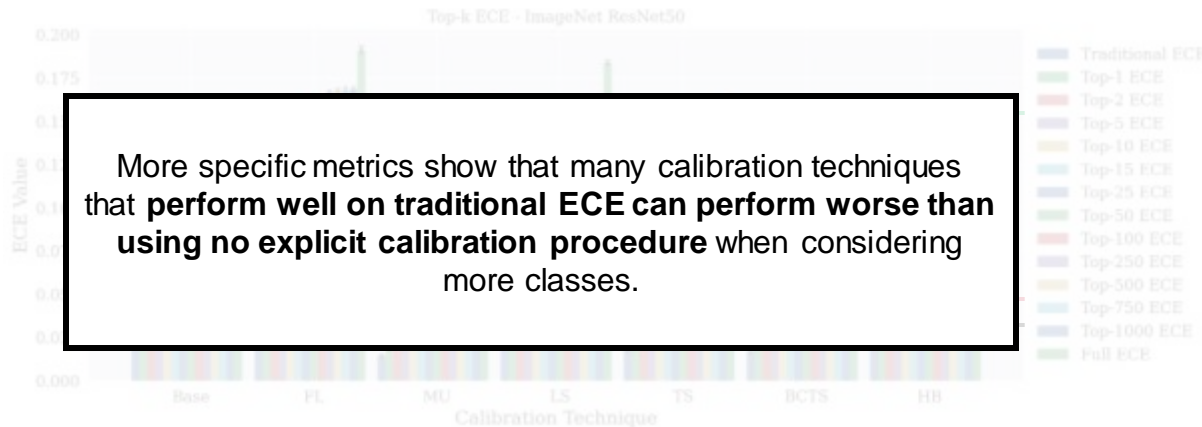
Question: How do these methods perform outside of the most confident class?

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes



Full ECE—Measures miscalibration across all classes

Overall Goal: Evaluate the state of the art in classifier calibration according to context focused metrics to observe how they perform in different definitions of reliability.

Experiment #1: Top- k

- Data Set: ImageNet
- Base Model: ResNet50

Question: How to these methods perform outside of the most confident class?

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes
4. How confidence will be shown to an end user

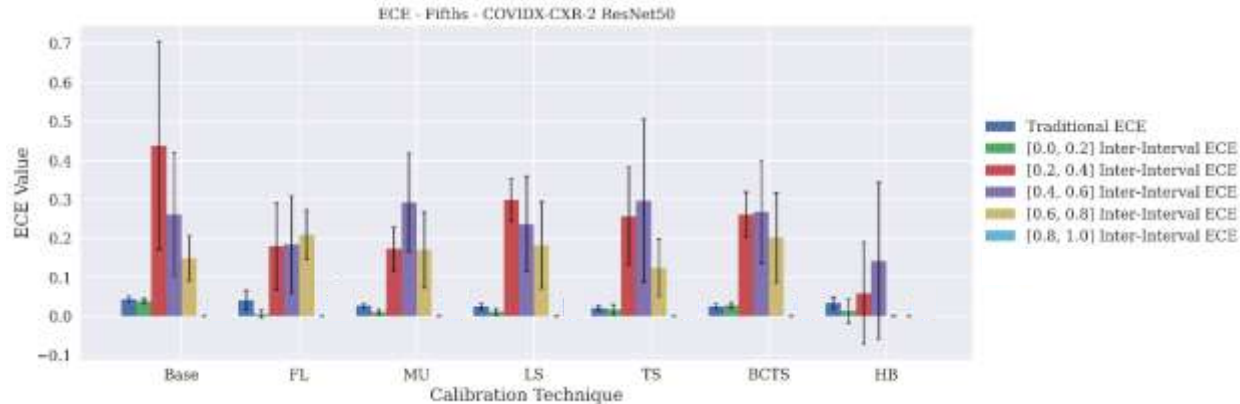
Experiment #2: Inter-Interval ECE

- Data Set: COVID-CRX-2
- Base Model: ResNet50

Assume:

Confidence will be displayed as to a clinician as one of five categories:
 [0.0,0.2] – Very low confidence of COVID
 [0.2,0.4] – Low confidence of COVID
 ...
 [0.8,1.0] – Very high confidence of COVID

How can we evaluate classifier calibration in this context?

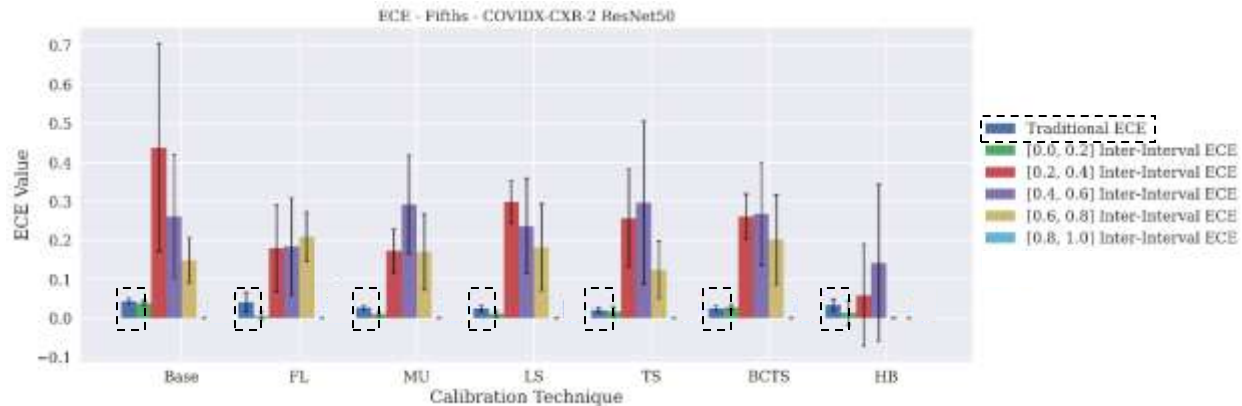


Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes
4. How confidence will be shown to an end user



Top-1 ECE – Measures miscalibration for most confident class

Experiment #2: Inter-Interval ECE

- Data Set: COVID-CRX-2
- Base Model: ResNet50

Assume:

Confidence will be displayed as to a clinician as one of five categories:
 [0.0,0.2] – Very low confidence of COVID
 [0.2,0.4] – Low confidence of COVID
 ...
 [0.8,1.0] – Very high confidence of COVID

How can we evaluate classifier calibration in this context?

Our Work: Context Focused Calibration Metrics

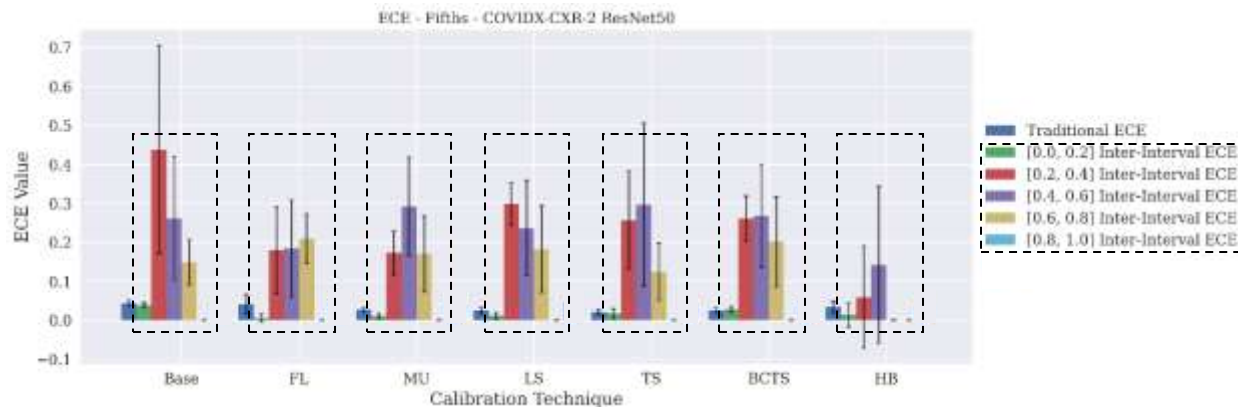
(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes
4. How confidence will be shown to an end user

Experiment #2: Inter-Interval ECE

- Data Set: COVID-CRX-2
- Base Model: ResNet50



Assume:

Confidence will be displayed as to a clinician as one of five categories:

[0.0,0.2] – Very low confidence of COVID

[0.2,0.4] – Low confidence of COVID

...

[0.8,1.0] – Very high confidence of COVID

How can we evaluate classifier calibration in this context?

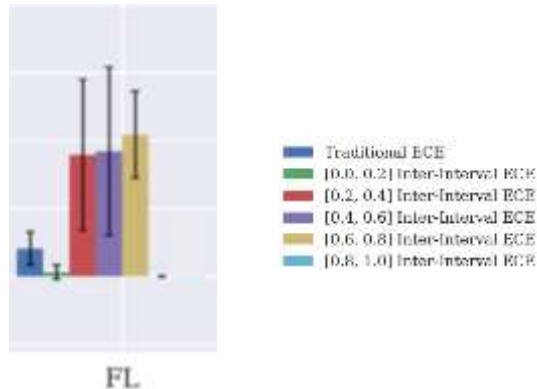
Inter-Interval ECE—Measures degree of miscalibration with respect to each category

Our Work: Context Focused Calibration Metrics

(Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes
4. How confidence will be shown to an end user



Experiment #2: Inter-Interval ECE

- Data Set: COVID-CRX-2
- Base Model: ResNet50

Assume:

Confidence will be displayed as to a clinician as one of five categories:

[0.0,0.2] – Very low confidence of COVID

[0.2,0.4] – Low confidence of COVID

...

[0.8,1.0] – Very high confidence of COVID

How can we evaluate classifier calibration in this context?

Inter-Interval ECE– Measures degree of miscalibration with respect to each category

Our Work: Context Focused Calibration Metrics (Kirchenbauer, Oaks, and Heim; 2021 *Under Review*)

Using a statistical framing of ECE, we developed a number of metrics that consider these factors:

1. Application-specific tradeoffs between classes (e.g., “Friendly” versus “Enemy” vehicles)
2. Specific instances of interest (e.g., Measuring calibration on instances with label “Enemy Vehicle”)
3. Subsets of the class probability space between most confident class and all classes
4. How confidence will be shown to an end user

Experiment #2: Inter-Interval ECE

- Data Set: COVID-CRX-2
- Base Model: ResNet50

Inter-Interval ECE enables evaluation of classifiers according to specified confidence categories that **reflect classifier usage**.

Assume:

Confidence will be displayed as to a clinician as one of five categories:
 [0.0,0.2] – Very low confidence of COVID
 [0.2,0.4] – Low confidence of COVID
 ...
 [0.8,1.0] – Very high confidence of COVID

How can we evaluate classifier calibration in this context?

Inter-Interval ECE– Measures degree of miscalibration with respect to each category

Team



Eric Heim
Senior ML Researcher
AI Division



John Kirchenbauer
Machine Learning Engineer
AI Division



Jon Helland
Machine Learning Researcher
AI Division



Jacob Oaks
Student Intern
AI Division



Aarti Singh
Associate Professor
Machine Learning Department



Zachary Lipton
Assistant Professor
Machine Learning Department

Final Thoughts

Machine-learned models are able to express **uncertainty** in their predictions that can lead to more informative, robust AI systems by

1. allowing humans to reason about when the model is likely to be incorrect
2. allowing components in a larger system to take different actions based on model confidence

In this project we research methods to evaluate, characterize, articulate and rectify uncertainty

Next steps:

- Develop a demonstration highlighting the utility of accurately expressing uncertainty.
- Create techniques to characterize the cause of uncertainty for a ML model.

For the audience: We are always looking for motivating real-world uses for our work. If you have a need for AI Systems that are able to express and reason under uncertainty, do not hesitate to reach out.

etheim@sei.cmu.edu