

# tinyML<sup>®</sup> Summit

*Miniature dreams can come true...*

**March 28-30, 2022 | San Francisco Bay Area**



[www.tinyML.org](http://www.tinyML.org)

# Accelerated Model Deployment with LEIP Recipes

Haya Sridharan,  
Product Manager

Latent AI  
March 2022

# Latent AI Overview

- We develop software to optimize AI models for low power and low latency.
- Founded Dec 2018 as SRI startup spin-off backed by DARPA technologies
- Team of proven startup-veterans, accomplished in AI. 25 FT staff (MS/PhD)
- Offices in Menlo Park, CA and Princeton, NJ

## Awards



<https://bit.ly/2EVXsBF>



Top 100 AI Startup  
2021



<https://tcn.ch/34juQee>  
<https://vimeo.com/460028482>



Top 60 Edge  
Computing  
Companies

## Investors



BLACKHORN  
VENTURES



AUTOTECH  
VENTURES



# An Edge ML Story



Solutions Steve  
(Links with customer)

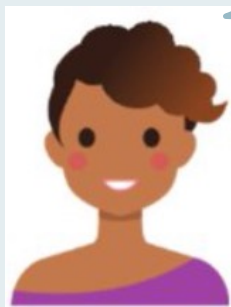
Let's detect  
potholes!

Ooh! I read a cool  
new paper, let me  
train a model :)



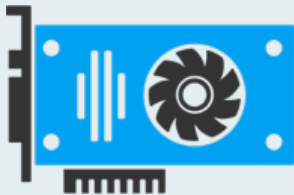
Data Scientist  
Daisy

Yikes!!



ML Engineer  
Evelyn

# Edge ML Challenges



Resource Constraints



Developer Frustration



Difficult to Iterate

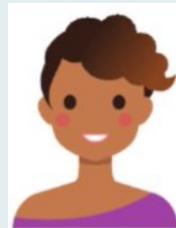
Challenges:

Pain points

- Inference Speed
- Model Size
- Memory
- Cost
- Power

- Hardware dependent Workflows
- Manual Optimization
- Poor Runtime Support

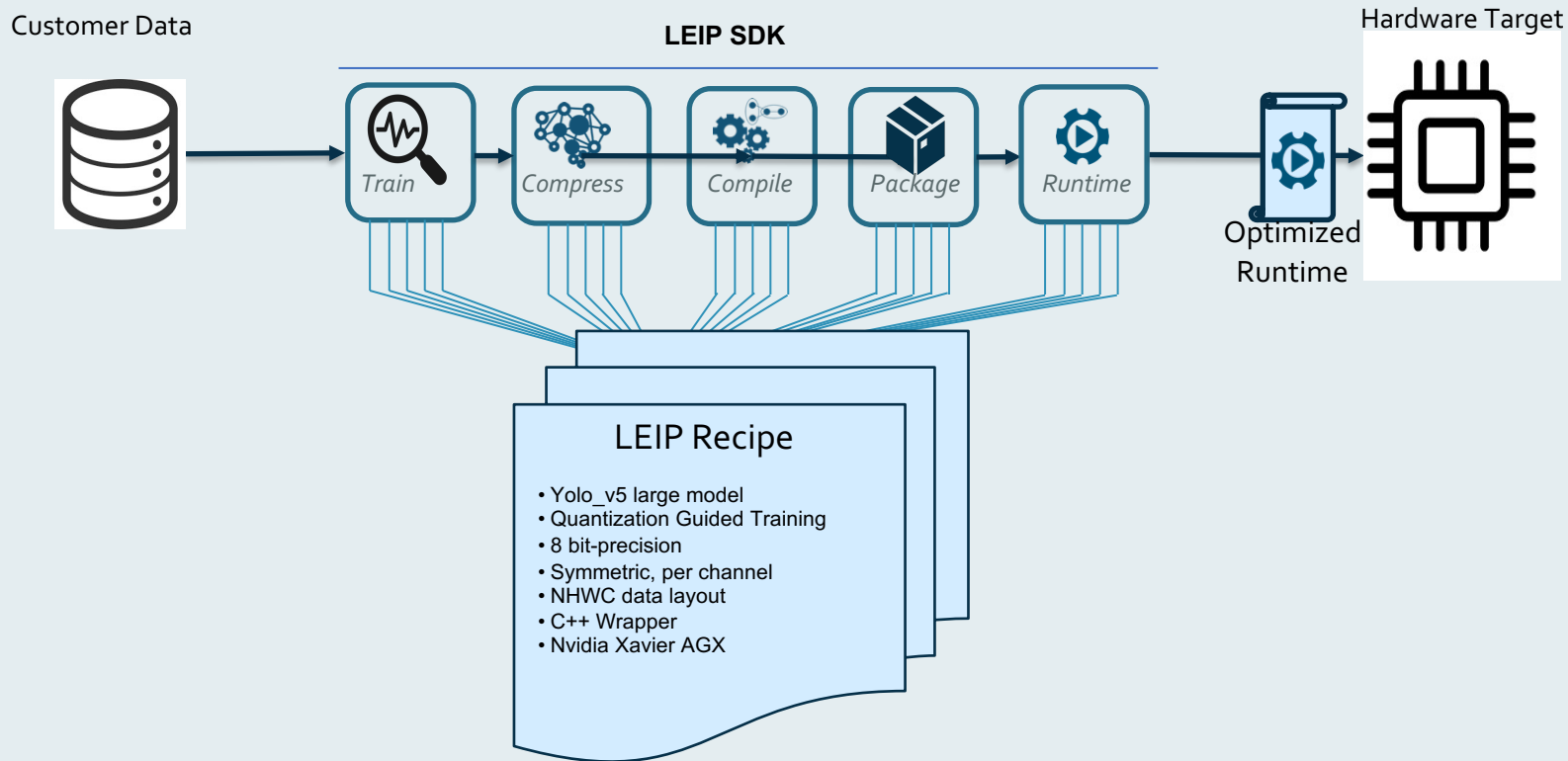
- Lack of Diagnostics
- Irreproducible Models
- Changing data landscape



ML Engineer  
Evelyn

Is there an easier way to do this?

# LEIP Recipes



# Recipe 1 - YoloV5 + AGX + COCO



**310%**

Inference Speed Improvement



**< 3%**

Accuracy Loss (mAP 0.5)



**84%**

Compression

- Pre-configured and optimized
- Demonstrates best-in-class performance
- Customers can BYOD

# Recipe Configuration

## Example: Train YOLOv5 on MS COCO

**af --config-dir=recipes --config-name=yolov5\_L\_RT**

recipes/yolov5\_L\_RT.yaml

```
defaults:
- hydra: defaults
- data: torchvision/coco-detection
- model: zhiqwang-yolov5-rt-stack
- trainer: ddp
- callbacks: defaults
- eval: coco
- predict: defaults
- vizdata: defaults
- export: defaults
- paths: defaults
- command: train
- override /hydra/job_logging: colorlog
- _self_

task:
width: 640
height: 640
batch_sizes: [8, 8]
num_workers: 4
moniker: recipe_${model.module.model_architecture}
```

## Data

af/configs/data/torchvision/coco-detection.yaml

```
nclasses: 80
module:
  _target_: af.core.data.modules.adaptermodule.AdapterDataModule
  batch_sizes: ${task.batch_sizes}
  num_workers: ${task.num_workers}
  adaptors: ${model.adaptors}
  train_transforms:
    _target_: af.core.data.augmentations.basic.resize
    width: ${task.width}
    height: ${task.height}
  valid_transforms:
    _target_: af.core.data.augmentations.basic.resize
    width: ${task.width}
    height: ${task.height}
  dataset_generator:
    _target_: af.core.data.sets.adaptors.torchvision-coco.MSCoco
    root_path: ${paths.cache_dir}
    task_family: "detection"
```

## Model

af/configs/model/zhiqwang-yolov5-rt-stack.yaml

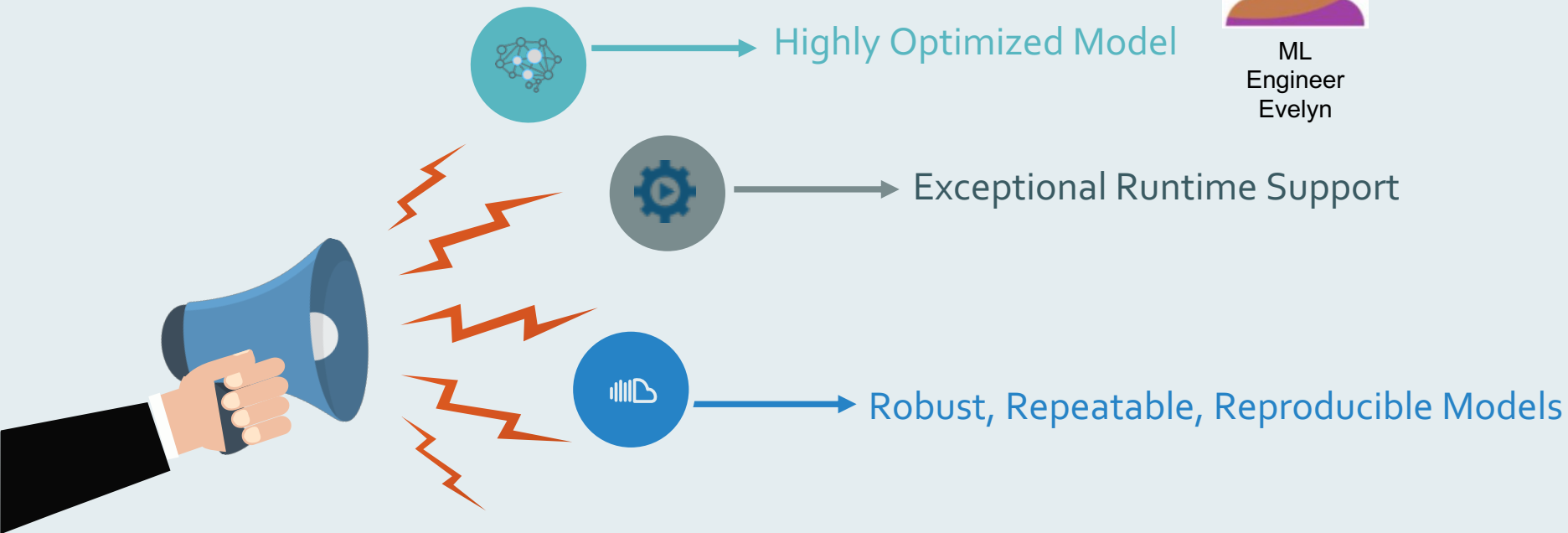
```
adaptors:
  uniform2trainvaltest:
    _target_: af.core.models.zoo.zhiqwang-yolov5-rt-stack.uniform2trainvaltest
  uniform2prediction:
    _target_: af.core.models.zoo.zhiqwang-yolov5-rt-stack.uniform2prediction
  prediction2uniform:
    _target_: af.core.models.zoo.zhiqwang-yolov5-rt-stack.prediction2uniform

trace_for_leip:
  _target_: af.core.models.zoo.zhiqwang-yolov5-rt-stack.trace_for_leip

module:
  _target_: af.core.models.zoo.zhiqwang-yolov5-rt-stack.YoloClassFactory
  patch_dir: ${paths.patch_dir}
  hash: 06022fd62dc247f1140e34bb0745673bce95ccad
  patch_file: zhiqwang-yolov5-rt-stack/20220225.patch
  model_architecture: yolov5s #yolov5s, yolov5m, yolov5l, yolov5x
  img_height: ${task.height}
  img_width: ${task.width}
  lr: 0.01
  score_thresh: 0.005 #0.005 is value tested for Coco
  num_classes: ${data.nclasses}
  pretrained: False # if you set pretrained: True, you must use the Coco dataset so num_classes match
```

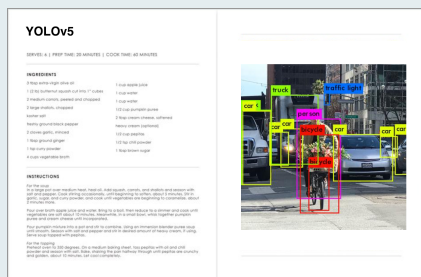


# Advantages of Recipes

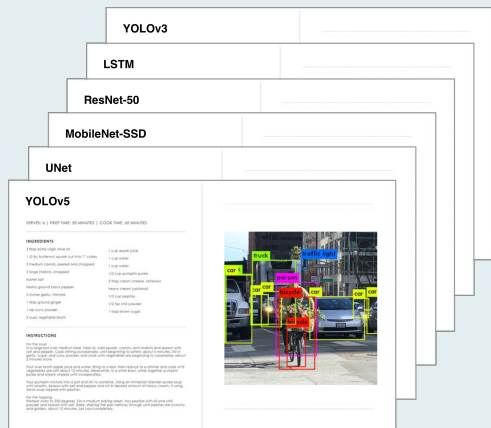


Recipes accelerate your time to market by 10x

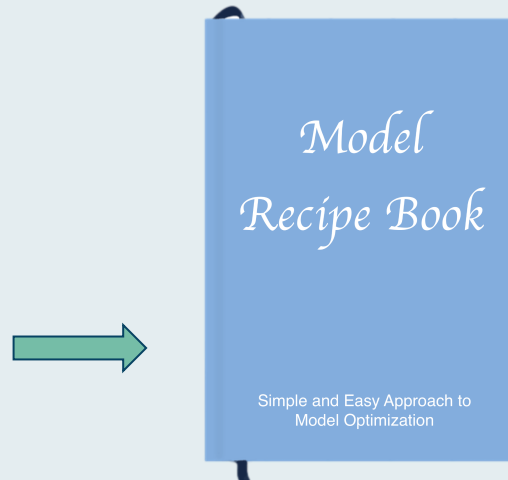
# Scaling Recipes



Optimize for a single model architecture and a hardware target



Optimize for multiple model architectures running across a heterogeneous compute environment



More models and more hardware targets with more optimization techniques including throttling and pruning

# Next Steps

If you want to -

- Reduce your Bill of Materials
- Fit multiple models into a single target hardware
- Accelerate your time-to-market for Edge AI Deployment

...we'd love to help! Please contact us at [info@latent.ai](mailto:info@latent.ai).

# Thank You

Jags Kandasamy, CEO  
jags@latent.ai  
+1 (404) 790-2048

Young Yoon, Products  
young@latent.ai  
+1 (510) 390-2041

Haya Sridharan, Products  
haya@latent.ai  
+1 (206) 330-7817





# Copyright Notice

This presentation in this publication was presented as a tinyML® Summit 2022. The content reflects the opinion of the author(s) and their respective companies. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

**[www.tinyml.org](http://www.tinyml.org)**