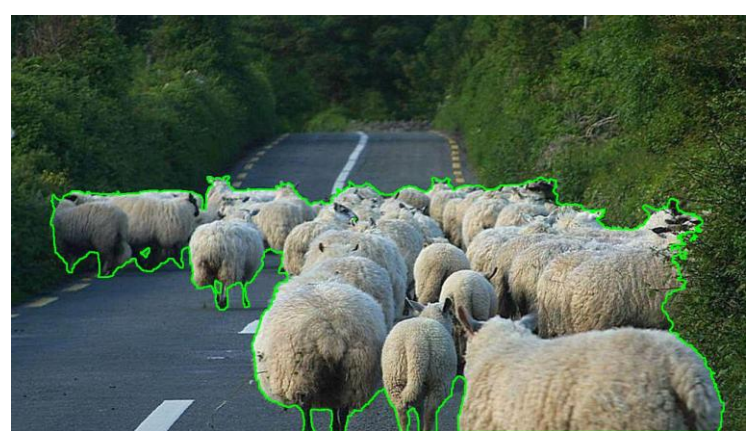




1. Hard Cases Detection in Motion Prediction by Vision-Language Foundation Models

Yi Yang, Qingwen Zhang, Kei Ikemura, Nazre Batool, John Folkesson

- Addressing hard cases is challenging!
 - Sparsity & high variability



Anomalous road users



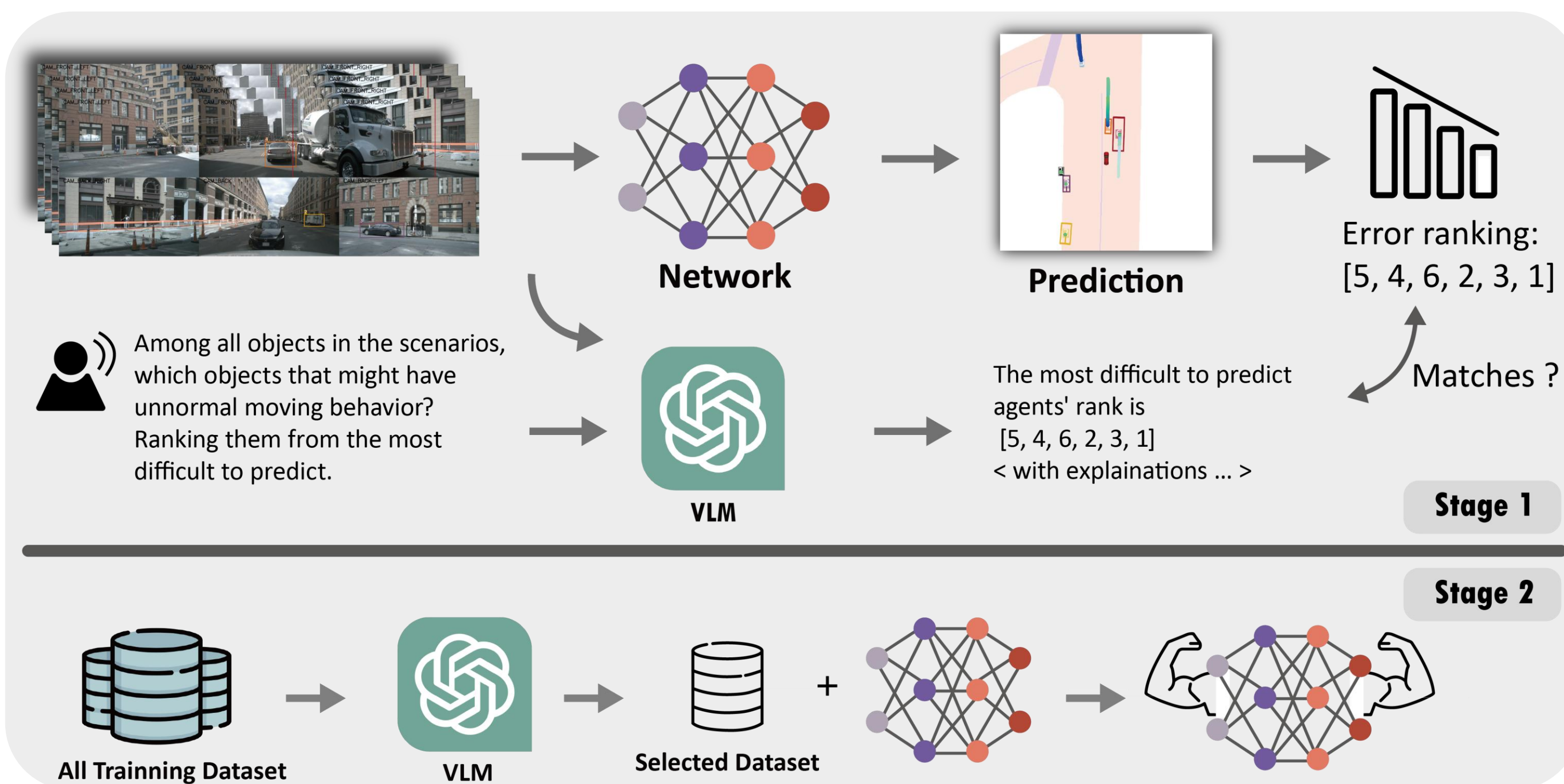
Extreme weather



Complex traffic

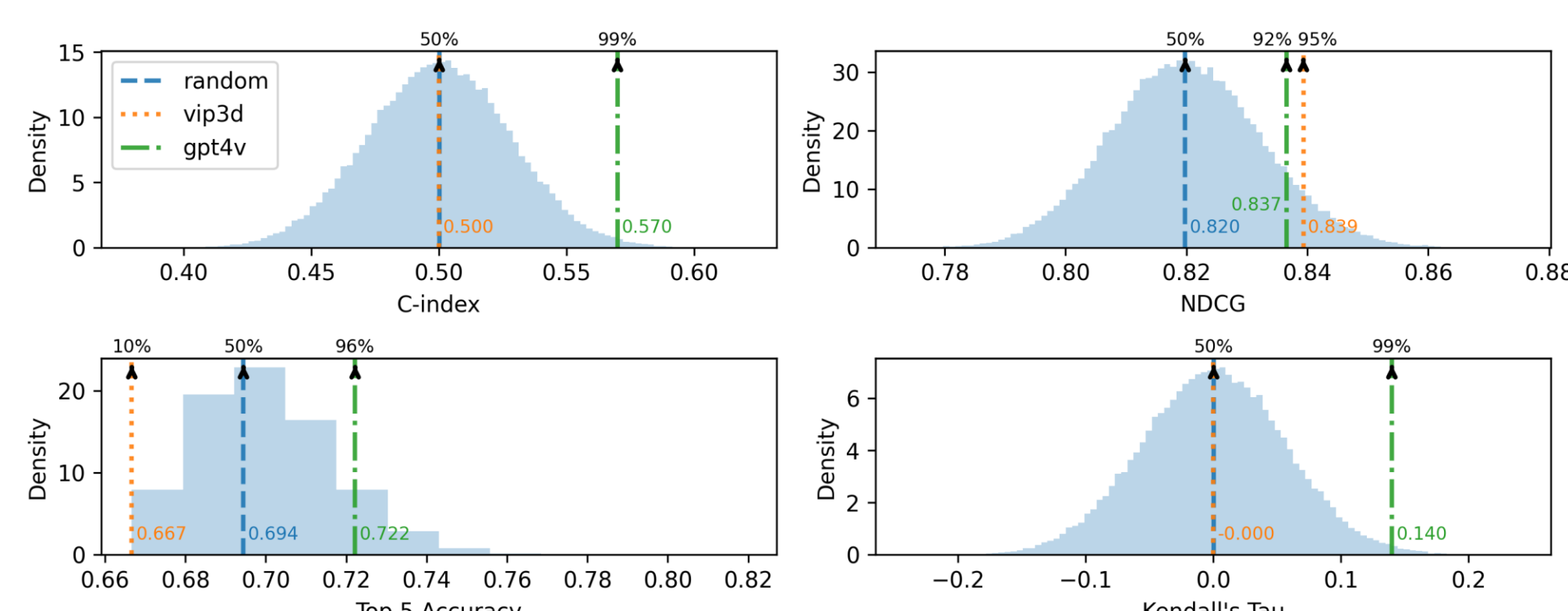
- Existing method:
 - collect more real-world data? -> *expensive!*
 - synthetic data?
 - generate with deep generative models conditioned on specific needs
 - manipulate the 3D reconstructed environment, like moving/adding road users
 - > *require much human intervention!*
 - Incremental learning? -> *dependence on the network training!*

Q: Is there a more explainable and independent method available?



Stage 1: Agent-level

- Verify** the ability of VLM to detect hard cases
- using existing motion prediction NN as ground truth



Stage 2: Scene-level

- Improve** training efficiency by training the network with a smaller subset of data selected by VLM.

	Class	Vehicle		Pedestrian	
	# Samples [Ratio%]	minADE	minFDE	minADE	minFDE
Whole	28130 [100]	0.71	1.02	0.82	1.11
Random	2000 [7.1]	0.80 \uparrow 13%	1.22 \uparrow 19%	0.93 \uparrow 13%	1.31 \uparrow 18%
	1000 [3.6]	0.82 \uparrow 16%	1.27 \uparrow 24%	0.92 \uparrow 12%	1.29 \uparrow 16%
	500 [1.8]	0.93 \uparrow 31%	1.45 \uparrow 42%	0.99 \uparrow 21%	1.42 \uparrow 28%
	200 [0.7]	0.97 \uparrow 36%	1.55 \uparrow 51%	1.03 \uparrow 25%	1.45 \uparrow 30%
GPT-4v	100 [0.4]	1.08 \uparrow 52%	1.73 \uparrow 69%	1.13 \uparrow 38%	1.63 \uparrow 47%
	200 [0.7]	0.93 \uparrow 31%	1.48 \uparrow 45%	1.03 \uparrow 25%	1.43 \uparrow 29%
	100 [0.4]	0.97 \uparrow 37%	1.56 \uparrow 52%	1.07 \uparrow 31%	1.56 \uparrow 40%

- 4 real examples of GPT4 outputs



"nighttime driving and wet road surfaces, which can affect visibility and vehicle behavior. The reflections and glare from the lights..."



"at an intersection... There is a large truck on the left that may obstruct the view and movement... increase the difficulty of prediction due to potential blind spots..."



"The intersection ahead adds complexity to the driving scenario, but overall traffic density is not high..."



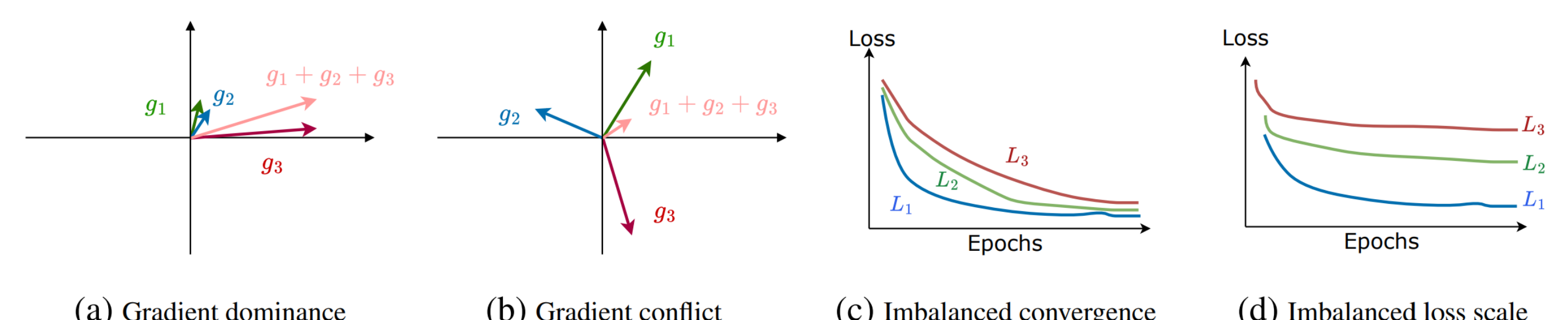
"The traffic situation appears to be straightforward with light traffic and clear road markings..."

2. AutoScale: Combining Multi-Task Optimization with Linear Scalarization

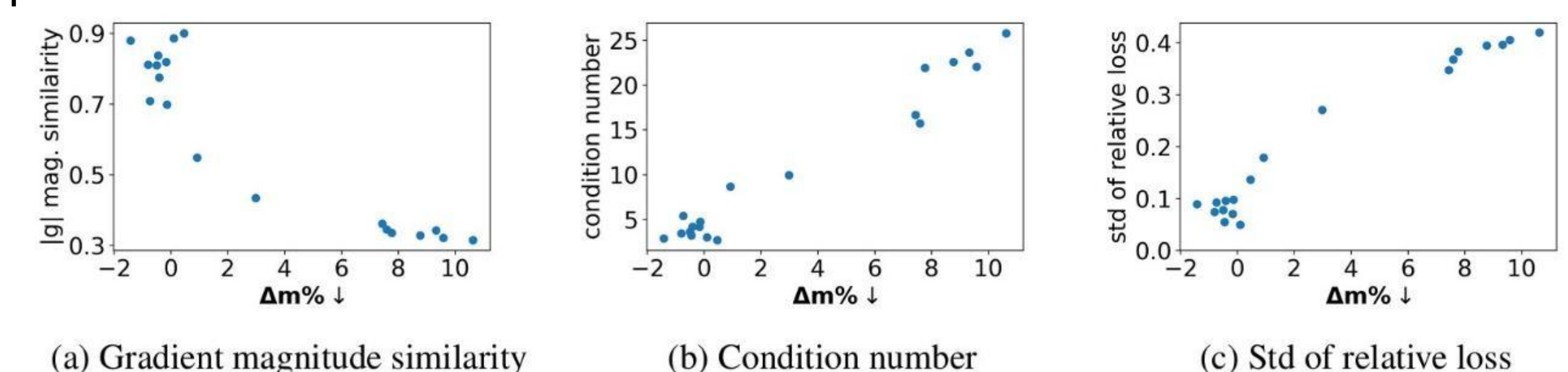
Yi Yang*, Kei Ikemura*, Qingwen Zhang, Ci Li, Nazre Batool, Sina Mansouri Sharif, John Folkesson

- There exists multiply MTL training issues!
 - Multi-objective learning: gradient & loss

$$\mathcal{L}(\theta; w) \equiv \sum_{i=1}^K w_i \mathcal{L}_i(\theta), \quad w > 0, \quad \sum_i w_i = 1.$$



- Surprisingly, we find some metrics serve as good indicators of performance!



- Given high performance -> optimal MTO metric values, we hypothesize that the reverse also holds: optimize MTO metric -> high performance. If so, we can localize the optimal task weights by optimizing the key metric value.
- We propose AutoScale, an efficient and practical two-stage pipeline that partitions a single training run into two phases.

$$w^* = \arg \min_w \mathbb{E}[F(w|\{\mathcal{G}\}, \{\mathcal{L}\})], \quad \text{s.t.} \quad \sum_{i=1}^K w_i = K,$$

Optimal weight set Given set of gradient and loss K loss terms in total

- Here we test three cost function considering gradient magnitude similarity, loss similarity, and condition number.