



MAP: Model Merging with Amortized Pareto Fronts Using Limited Computation

Lu Li^{1*}, Tianyu Zhang^{2,3*}, Zhiqi Bu^{4*}, Suyuchen Wang⁵, Huan He¹, Jie Fu⁶, Jiang Bian⁷, Yonghui Wu⁷, Yong Chen¹, Yoshua Bengio²





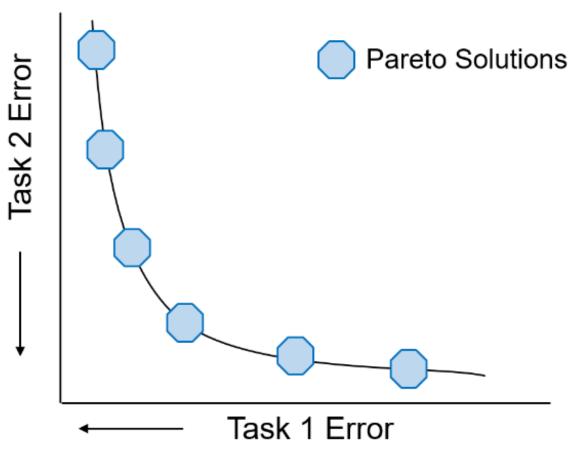
Results

Motivation

- Task-vector based model merging combines multiple single-task models into a single merged model that is capable of multi-task learning, without the need for sharing training data or any further training.
- The scaling coefficient for each task determines the relative performance of the merged model on that task

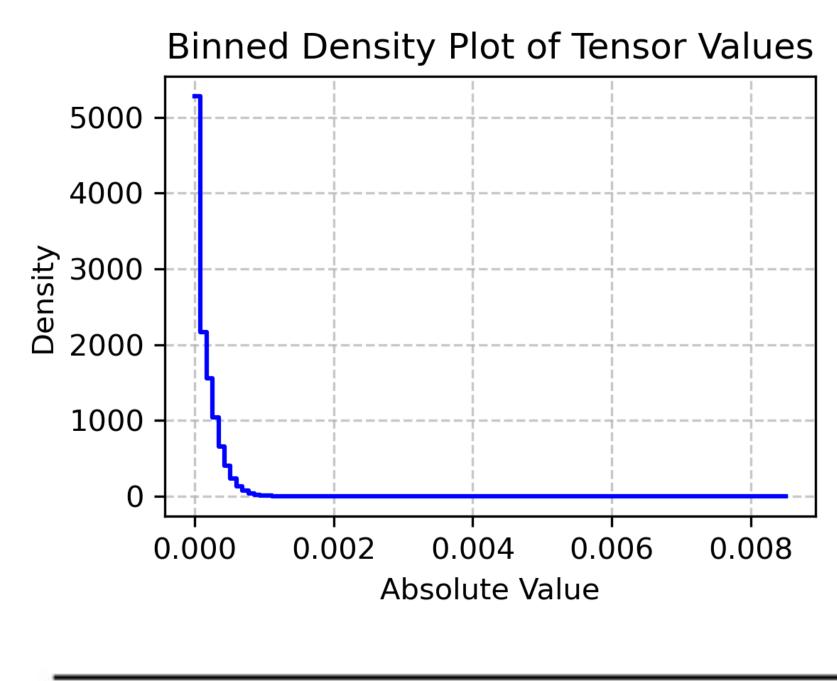
$$\boldsymbol{\theta}_m = \boldsymbol{\theta}_{pre} + \sum_i \boldsymbol{c_i} \boldsymbol{v_i}$$

- However, practitioners might have different preferences for the tasks, leading to trade-offs in how to merge the models. A set of Pareto optimal solutions is preferable.
- However, computing a Pareto front using conventional method involves many inferences and time-consuming



Key observation

• Fine-tuned models tend to converge near the pretrained model in parameter space.

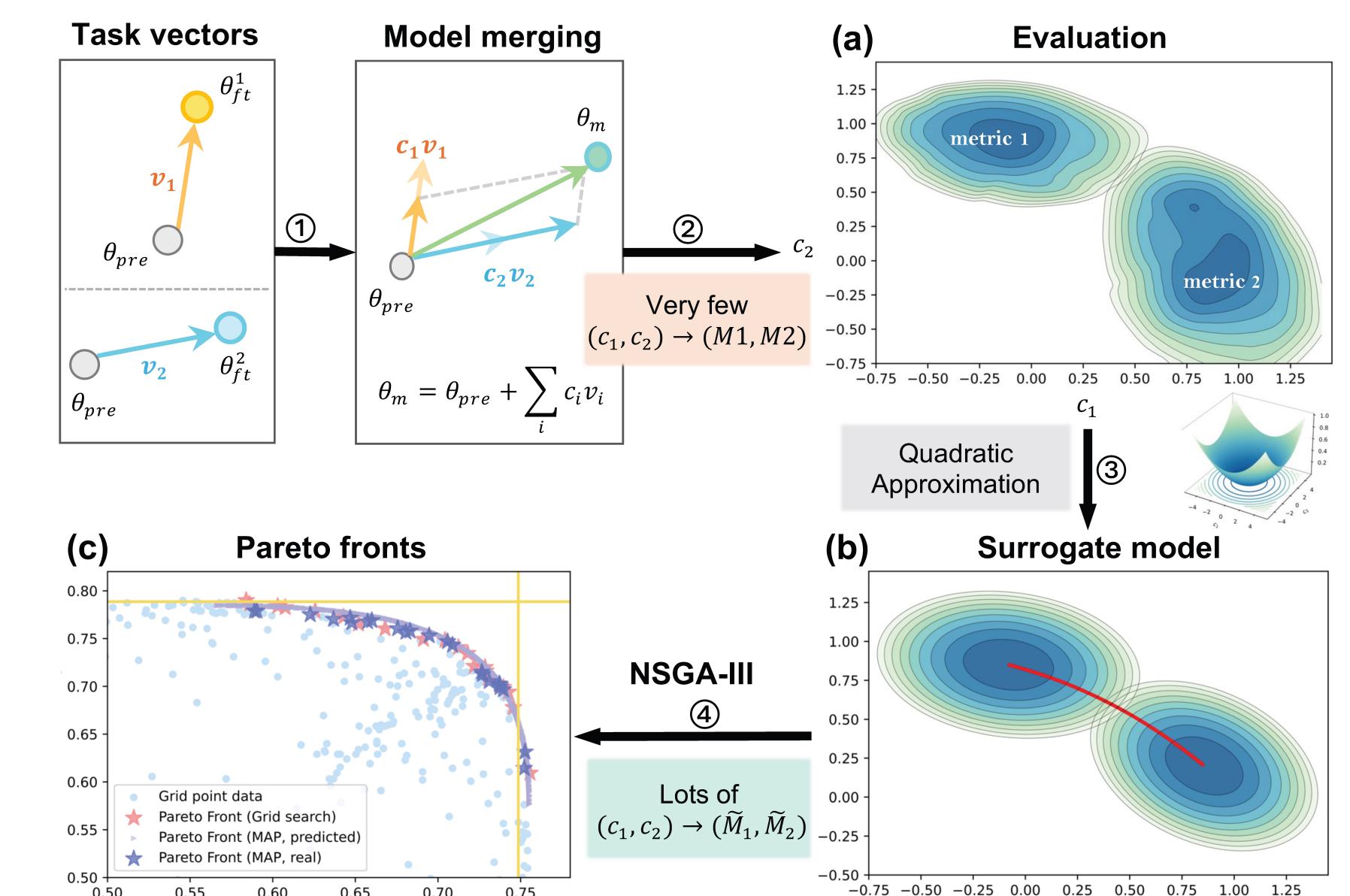


This motivates us to use second-order Taylor expansion to approximate the evaluation metrics for each task given a merged model

Metric	SUN397	Cars	DTD	SVHN
$ oldsymbol{ heta}_{pre} _1$	1,270,487	1,270,487	1,270,487	1,270,487
$ oldsymbol{ heta}_{pre} _1 \ \mathbf{v}_n _1$	21,055	20,127	13,621	19,349
$ \mathbf{v}_n _1/ \boldsymbol{\theta}_{pre} _1(\%)$	1.66%	1.58%	1.07%	1.52%

Method

 Task-vector based model merging combines multiple
 We propose MAP, a computationally efficient method to identify single-task models into a single merged model that is Pareto fronts, allowing trade-offs without requiring additional training.

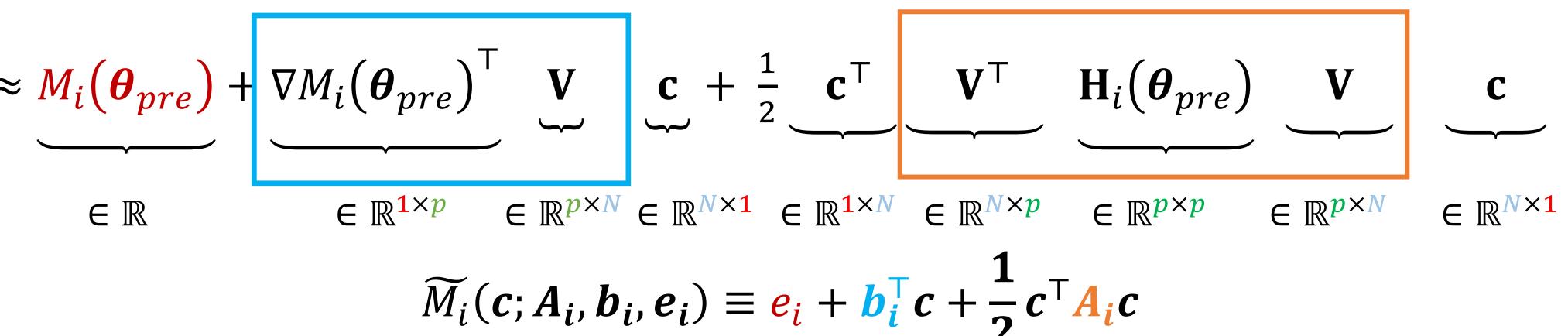


Let $M_i(\boldsymbol{\theta}_m(\boldsymbol{c})) = M_i(\boldsymbol{\theta}_{pre} + \sum_i \boldsymbol{c}_i \boldsymbol{v}_i)$ denote the evaluation metrics on task i

•
$$M_i(c) = M_i(\boldsymbol{\theta}_m(\mathbf{c}))$$

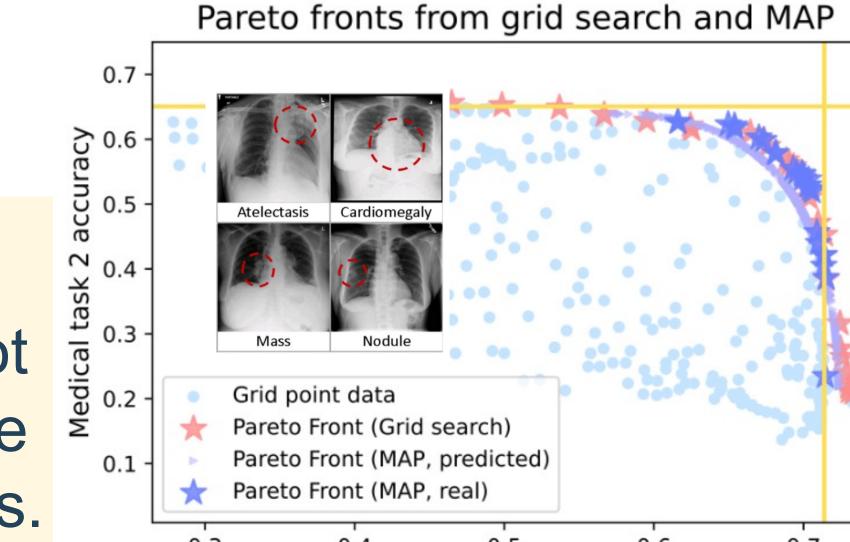
$$= M_i(\boldsymbol{\theta}_{pre}) + \nabla M_i(\boldsymbol{\theta}_{pre})^{\mathsf{T}} (\boldsymbol{\theta}_m(\mathbf{c}) - \boldsymbol{\theta}_{pre}) +$$

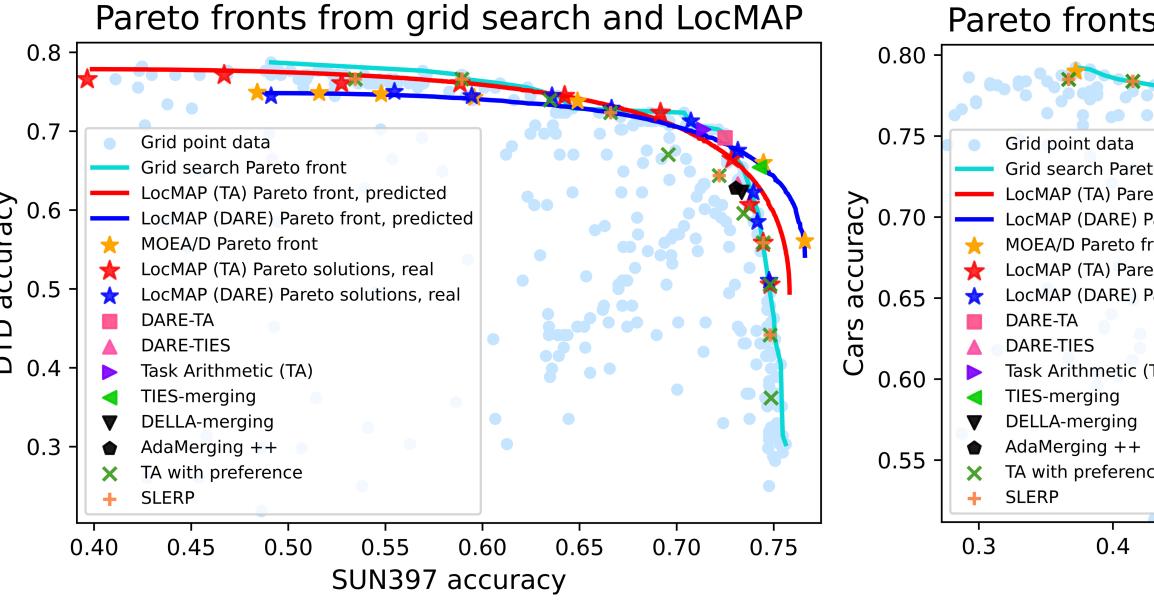
$$= \frac{1}{2} (\boldsymbol{\theta}_m(\mathbf{c}) - \boldsymbol{\theta}_{pre})^{\mathsf{T}} \mathbf{H}_n(\boldsymbol{\theta}_{pre}) (\boldsymbol{\theta}_m(\mathbf{c}) - \boldsymbol{\theta}_{pre}) + R_i(\boldsymbol{\theta}_m(\mathbf{c}) - \boldsymbol{\theta}_{pre})$$

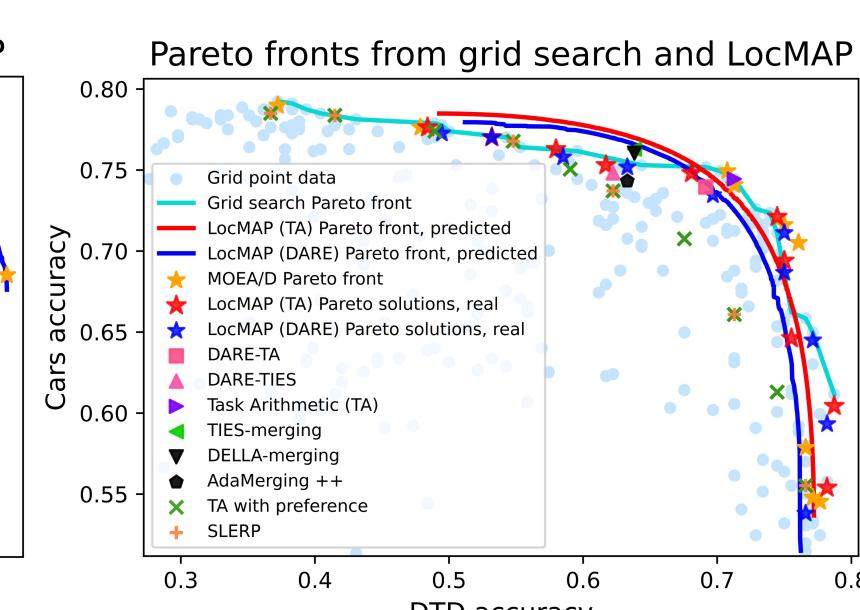


- Because of symmetry of A, the number of variables need to estimate is $\frac{(N+1)(N+2)}{2}$
- Since we already have the c and $M_i(c)$, it can be formulated as a linear regression, with $X=\left(c_1^2,c_2^2,\ldots,c_T^2,c_1c_2,c_1c_3,\ldots,c_{T-1}c_T,c_1,c_2,\ldots,c_T,1\right) \text{ and } y=M_i\left(\theta_m(c)\right)$
- Closed form solution: $(X^TX)^{-1}X^Ty$

MAP is able to find diverse Pareto front, not captured by other single merged model baselines.







Medical task 1 accuracy

For higher dimensions, where we cannot visualize the Pareto front, the **win rate** is used to measure how often MAP outperforms other methods in terms of Pareto front solutions across tasks.

Win Rate =
$$\frac{1}{K^2 N} \sum_{i=1}^K \sum_{j=1}^K \sum_{n=1}^N \mathbb{I}\left[M_n(\theta(\mathbf{c}_i^{\text{LocMAP}})) > M_n(\theta(\mathbf{c}_j^{\text{baseline}}))\right]$$

N # c (direct search) # c per dim # c (MAP) Win rate (MAP)

2 3	200 300	14.14 6.69	30 50	$49.81\% (\pm 0.30)$ $46.90\% (\pm 0.71)$	$0.953 (\pm 0.018)$ $0.980 (\pm 0.003)$
$\frac{-}{N}$	# c (direct search)	# c per dim	# c (MAP)	Win rate (MAP)	R^2 (MAP)
4	300	4.16	60	50.67% (±2.44)	$0.984 (\pm 0.004)$
5	500	3.47	85	$53.00\% (\pm 1.88)$	$0.941 (\pm 0.001)$
6	500	2.82	100	$60.71\%~(\pm 1.34)$	$0.941\ (\pm0.030)$
7	1000	2.68	140	$63.42\% \ (\pm 1.91)$	$0.891 (\pm 0.024)$
8	1000	2.37	250	$65.58\% \ (\pm 0.94)$	$0.868~(\pm 0.028)$







 R^2 (MAP)