# Rewarded soups:

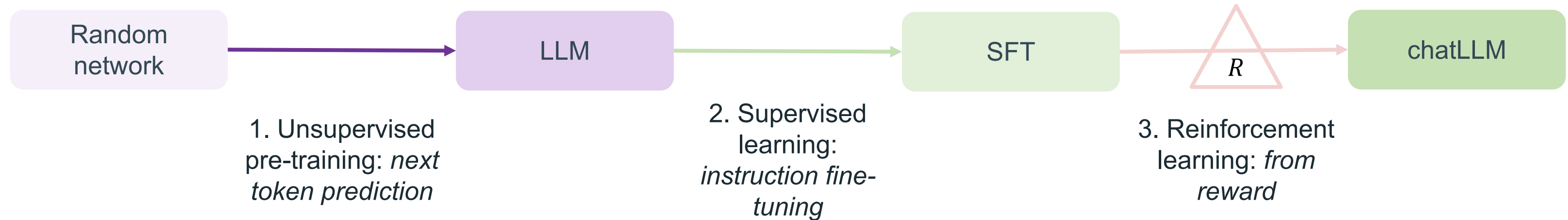## towards Pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards.

**Alexandre Ramé**, Guillaume Couairon, Corentin Dancette, Jean-Baptiste Gaya, Mustafa Shukor, Laure Soulier, Matthieu Cord.
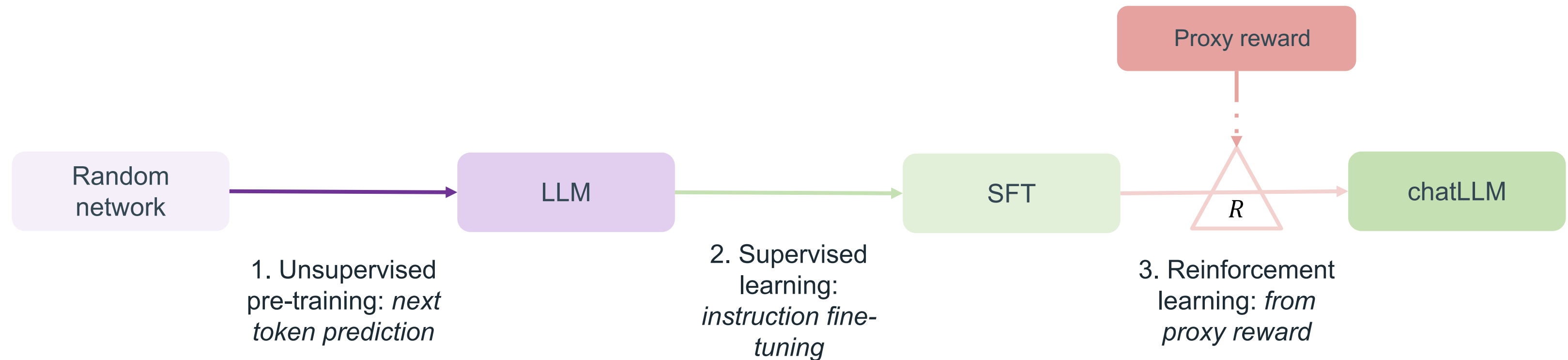
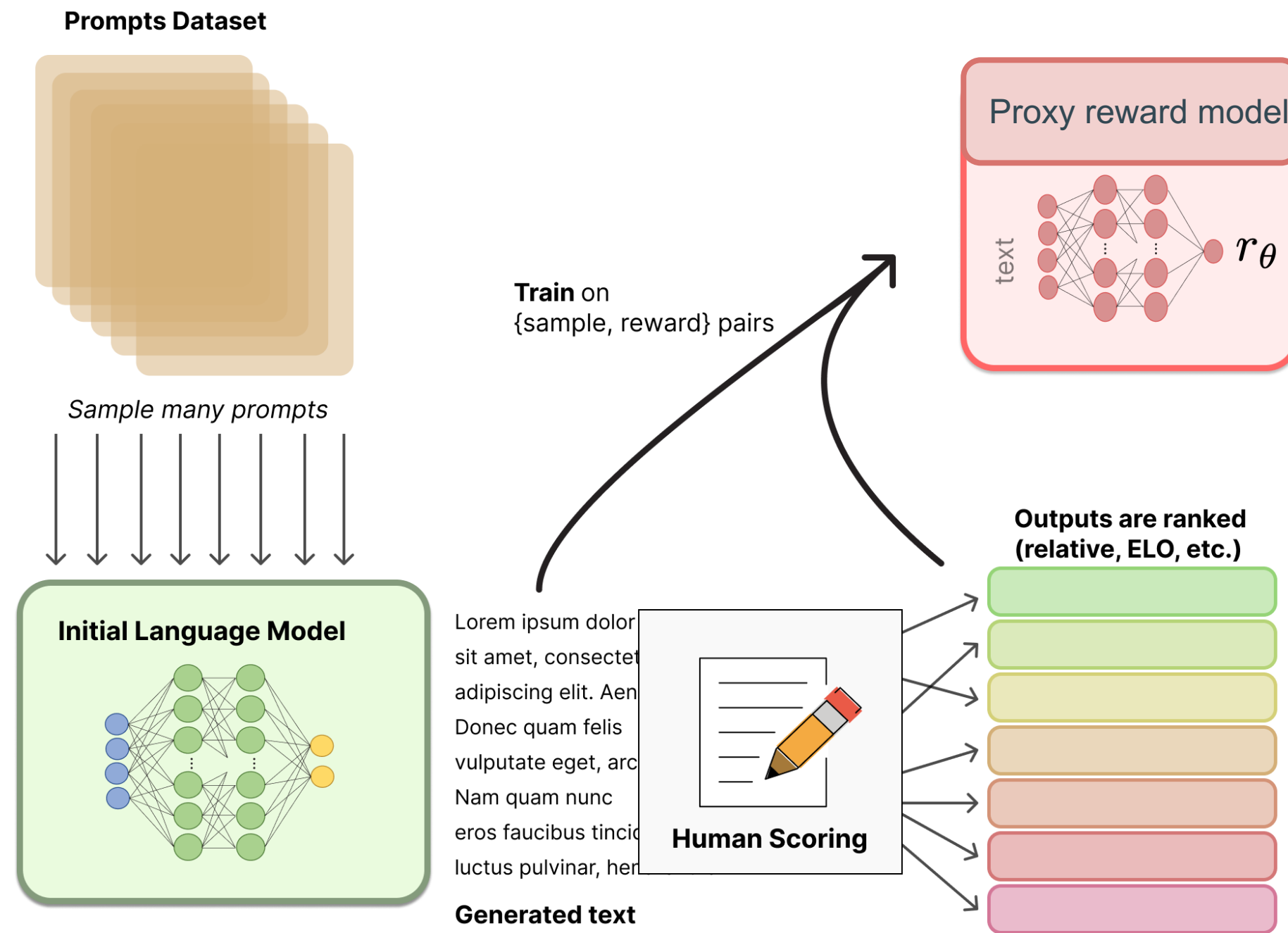NeurIPS 2023

# LLM training in 3 steps



Random network → LLM → SFT → R → chatLLM

1. Unsupervised pre-training: *next token prediction*

2. Supervised learning: *instruction fine-tuning*

3. Reinforcement learning: *from reward*

[Stiennon2020] Learning to summarize from human feedback. NeurIPS.
[Ouyang2022] Training language models to follow instructions with human feedback. NeurIPS.

# Need for a proxy reward in the RL step



**Problem**: the true reward is not available.
**Consequence**: proxy reward.
**Challenge:** reward misspecification
(when the training reward is not a good proxy).

[Stiennon2020] Learning to summarize from human feedback. NeurIPS.
[Ouyang2022] Training language models to follow instructions with human feedback. NeurIPS.

# Reward model from human feedback for RLHF



Prompts Dataset

Proxy reward model

text $r_\theta$

Train on
{sample, reward} pairs

Sample many prompts

Initial Language Model

Lorem ipsum dolor
sit amet, consectet
adipiscing elit. Aen
Donec quam felis
vulputate eget, arc
Nam quam nunc
eros faucibus tincic
luctus pulvinar, her

Generated text

Human Scoring

Outputs are ranked
(relative, ELO, etc.)

[Christiano2017] Deep reinforcement learning from human preferences. NeurIPS.

# Diversity of opinions

Consistency issue: only ≈65% agreement across labellers.

Indeed, human opinions are diverse (and subjective):

- **Politics**: democrat vs republican?
- **Uncertain situations**: economic strategy for climate change?
- **Aesthetic**: beautiful vs ugly?

More generally, different expectations from machines:

- **Safety**: helpfulness vs harmlessness?
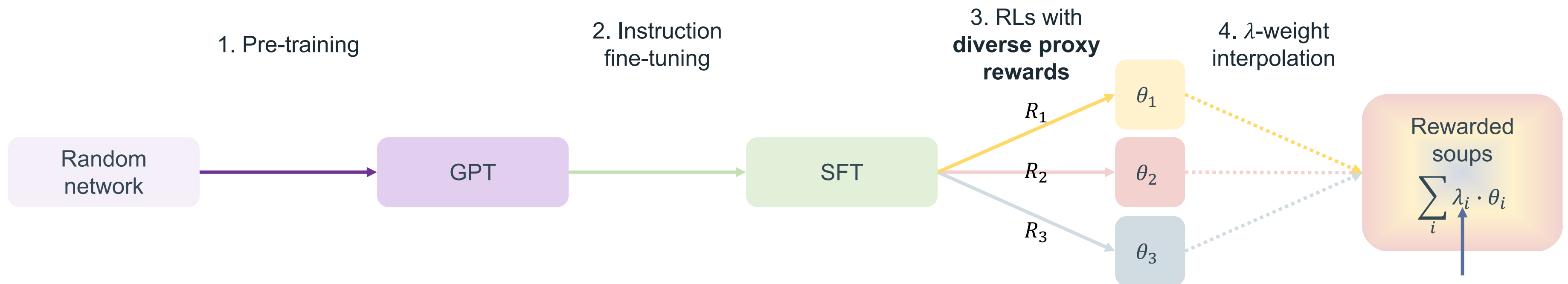- **Summarization**: complete or factual ?

Diversity of people and applications ⇒ which one should we optimize for?

"Human aligned artificial intelligence is a multi-objective problem" [Vamplew2018].

Move from a single-policy towards a **multi-policy paradigm** to embrace diversity.

# Rewarded soups: towards Pareto-optimal alignment by interpolating weights fine-tuned on diverse rewards
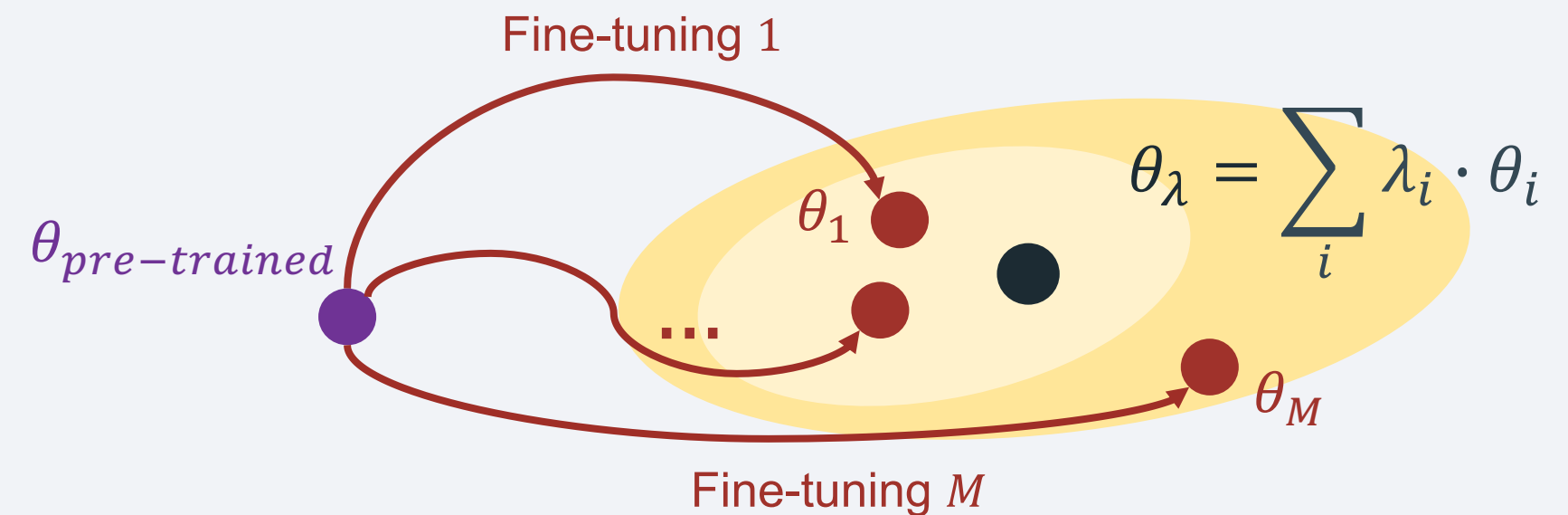
1. Pre-training

2. Instruction fine-tuning

3. RLs with **diverse proxy rewards**

4. $\lambda$-weight interpolation

Random network

GPT

SFT

$R_1$ $\theta_1$

$R_2$ $\theta_2$

$R_3$ $\theta_3$

Rewarded soups $\sum_i \lambda_i \cdot \theta_i$

5. The $\{\lambda_i\}_i$ selected by the users according to their preferences.

Rewarded soup:

1. From a shared pre-trained foundation model,
2. Fine-tuned to follow instructions,
3. Launch one RL fine-tuning for each reward, each representing an opinion,
4. Interpolate the weights expert on diverse rewards,
5. Reveal the front of solutions (and select one interpolating coefficient).

# Weight interpolation relies on linear mode connectivity



When fine-tuned from a shared **pre-trained** model, weights remain **linearly connected** and thus can be interpolated despite the non-linearities in the architecture.
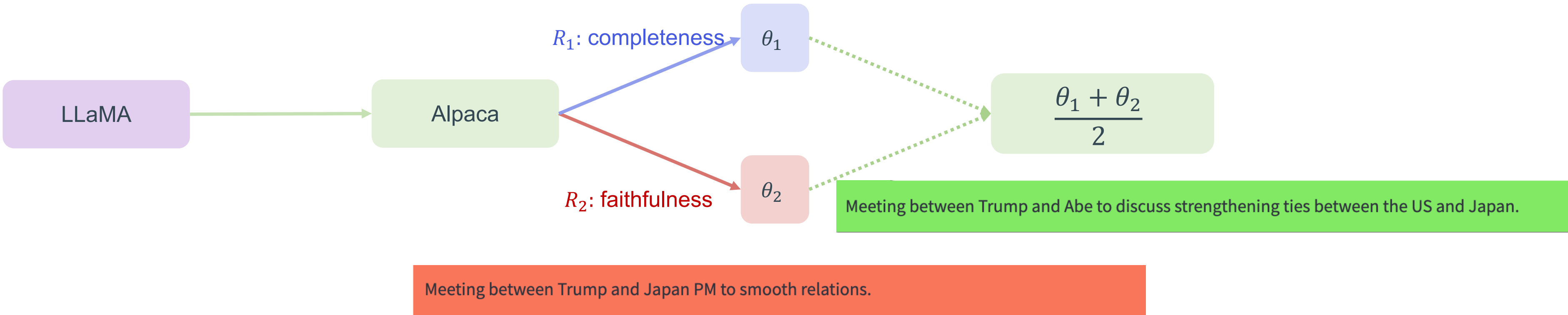
[Wortsman2022] Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. ICML.
[Rame2023] DiWA: diverse weight averaging for out-of-distribution generalization. NeurIPS.

# Summarization with diverse reward models

# Summarization with diverse reward models

# Pareto-optimal alignment across rewards



**Rewarded soups**

Interpolate the weights a posteriori:
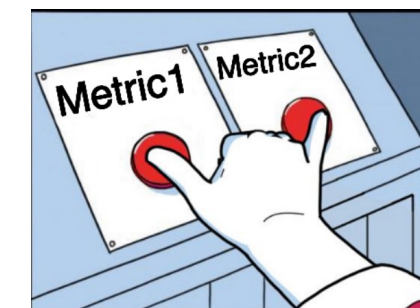
$$\sum_i \lambda_i \cdot \theta_i$$

**Multi-objective: MORL**

Interpolate the rewards a priori:

$$\sum_i \mu_i \cdot R_i$$

**Issue**: cost, as preference variations result in different solutions, requiring a high level of granularity.

In the plot:
- $R_2$ rewarded $(\mu = \lambda = 1)$
- MORL for $\mu = 0.5$
- RS for $\lambda = 0.5$
- RS front: $\{(1-\lambda) \cdot \theta_1 + \lambda \cdot \theta_2\}_\lambda$
- MORL front: $\{\theta_{(1-\mu)\times R_1 + \mu \times R_2}\}_\mu$
- LLaMA init
- $R_1$ rewarded $(\mu = \lambda = 0)$

In the paper, we theoretically prove the (approximated) Pareto-optimality of rewarded soups for quadratic rewards.

# We apply rewarded soup in multiple standard learning tasks:

## 1. Text

- Summarization (news, reddit).
- Movie review generation.
- Q&As of technical questions.
- Conversational assistant.

## 2. Multimodal: text and image

- Image captioning.
- Image generation with diffusion.
- Visual grounding.
- Visual question answering.

## 3. Continuous control

- Locomotion.

# Benefits from rewarded soups

## 1. Efficiency

- 1 single fine-tuning per reward.
- Parralelization.
- No inference overhead.

## 2. Transparence

- Support decision-making.
- Facilitate regulation by an (external) non-technical committee.

## 3. Updatable

- Easily update the $\lambda$.
- Easily add new reward.
- Iterative development process.

## 4. Fairness

- Value pluralism.
- Under-represented groups.
- Less ideological hegemony.
- Federated learning setups?.

Conclusion

- Human-aligned AI as a multi-objective problem.
- Weight interpolation as a practical pareto-optimal solution.
- Code: https://github.com/alexrame/rewardedsoups