



# Different Designs For LLM KD Loss

Alephia 25/5/5

# KLD



$$KL(p(X), q_{\theta}(X)) = E_{x \sim p(X)} \left[ \log \frac{p(x)}{q_{\theta}(x)} \right]$$

$$\begin{aligned} \operatorname{argmin}_{\theta} KL(p(X), q_{\theta}(X)) &= \operatorname{argmin}_{\theta} E_{x \sim p(X)} [-\log q_{\theta}(x)] \\ &= \operatorname{argmax}_{\theta} E_{x \sim p(X)} [\log q_{\theta}(x)] \\ &\approx \operatorname{argmax}_{\theta} \sum_x \log q_{\theta}(x) \\ &= \operatorname{argmax}_{\theta} \prod_x q_{\theta}(x) \end{aligned}$$

最小化KLD(p,q)等价于最小化CE(p,q)等价于最大化似然函数

# RKLD



最小化RKLD(p,q)等价于最小化CE(q,p)-H(q)

$$\begin{aligned} RKL(p(X), q_{\theta}(X)) &= KL(q_{\theta}(X), p(X)) \\ &= E_{x \sim q_{\theta}(X)} \left[ \log \frac{q_{\theta}(x)}{p(x)} \right] \\ &= E_{x \sim q_{\theta}(X)} [-\log p(x)] - H(q_{\theta}(x)) \end{aligned}$$

# FKLD: MEAN-SEEKING BEHAVIOUR

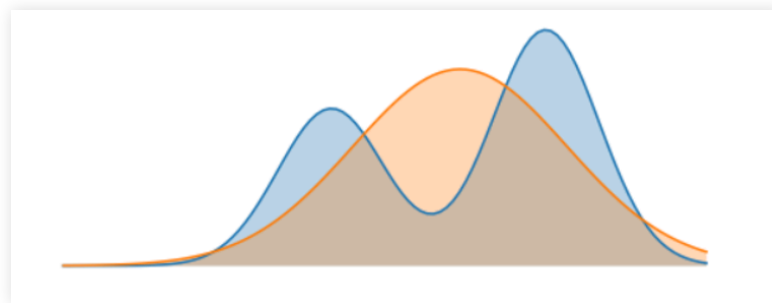


$$KL(p(X), q_{\theta}(X)) = E_{x \sim p(X)} [-\log q_{\theta}(x)] - H(p(x))$$

Zero Avoiding

$$\exists(x, y) \text{ s.t. } p(y|x) \gg 0, q_{\theta}(y|x) \approx 0 \rightarrow KL(p, q_{\theta}) = \inf$$

- p中高概率的地方，q也必须高，需要涵盖所有高概率区域
- q中高概率的地方，p不必高
- FKLD倾向于拟合多个峰



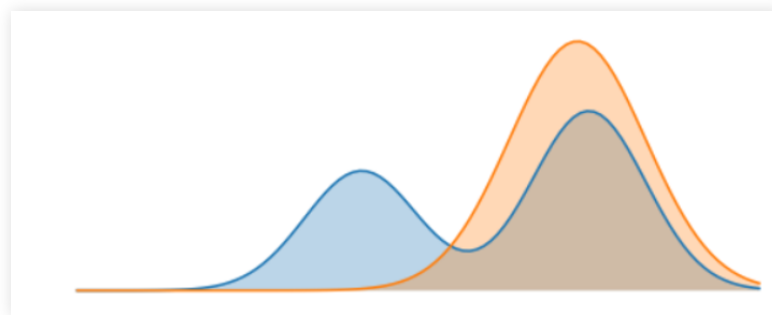
# RKLD: MODE-SEEKING BEHAVIOUR



$$RKL(p(X), q_{\theta}(X)) = E_{x \sim q_{\theta}(X)} [-\log p(x)] - H(q_{\theta}(x))$$

$$\exists(x, y) \text{ s.t. } q_{\theta}(y|x) \gg 0, p(y|x) \approx 0 \rightarrow KL(q_{\theta}, p) = \inf$$

- q中高概率的地方，p也必须高，q中低概率的地方，p也应该较小
- p中高概率的地方，q不必高
- RKLD倾向于拟合一个峰



# RKLD IN LLM KD



KLD下，学生在教师分布的viod region会高估，进而带来麻烦。这一问题在RKLD下有所缓解

条件： 1 教师服从混合Gaussian分布，学生服从Gaussian分布 2 两个分布都是连续的

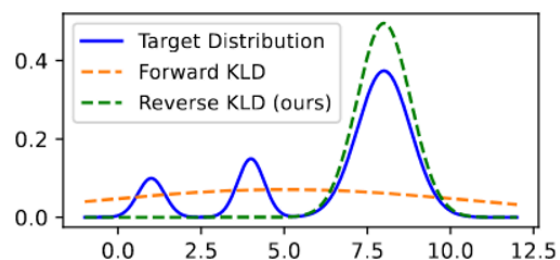


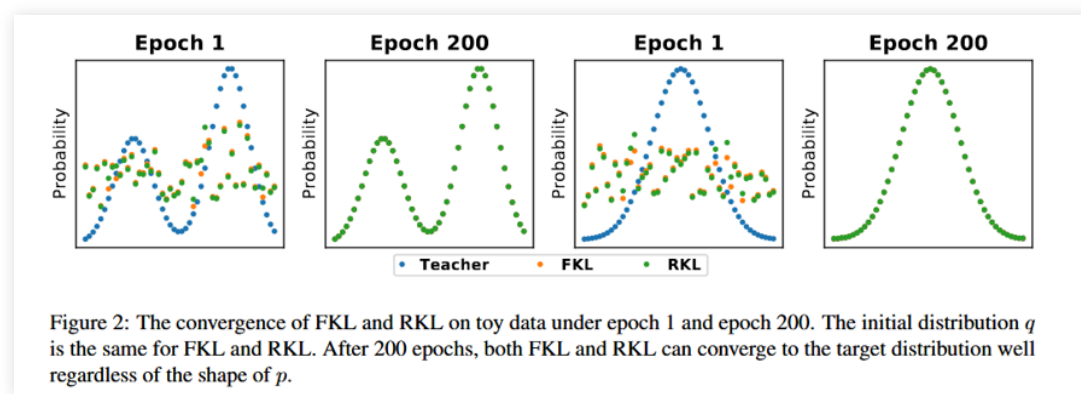
Figure 2: We fit a Gaussian mixture distribution with a single Gaussian distribution using *forward* KLD and *reverse* KLD.

Gu, Y., Dong, L., Wei, MiniLLM: Knowledge Distillation of Large Language Models. In ICLR,24

# DOES RKLD REALLY HELPS IN LLM KD?

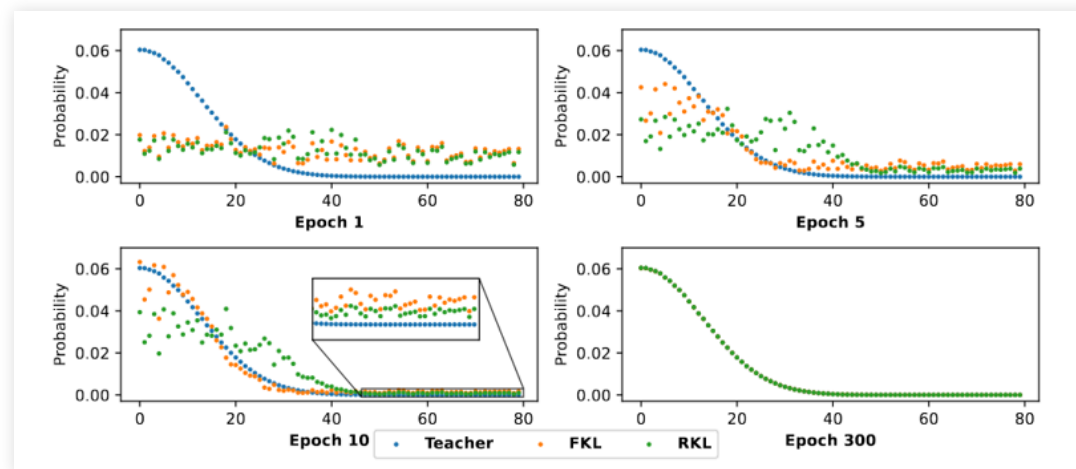
1. 教师，学生输出经过softmax之后不一定满足Gaussian分布
2. logits分布是离散的

事实上非Gaussian+离散情况下，充分训练后，两种loss训练下都会得到同一个拟合结果



Wu, T., Tao, Rethinking Kullback-Leibler Divergence in Knowledge Distillation for Large Language Models. In COLING,25

# COMBINE RKLD WITH FKLD



LLM KD中，所谓mean-seeking和mode-seeking可能并不存在，取而代之的是：FKLD倾向于先拟合分布头部，RKLD倾向于先拟合分布尾部

最终solution:  $AKL(p, q_\theta) = \alpha_1 FKL(p, q_\theta) + \alpha_2 RKL(p, q_\theta)$