

数理科学チームでは、幅広い純粋数学・理論物理の研究者の力を借りて、人工知能・機械学習分野における様々な数学的課題に組織的に取り組んでいる。このポスターでは、数理科学チームのE. M. Kiralと近似ベイズ推論チームのT. MöllenhoffとE. Khanの共同研究を解説する。

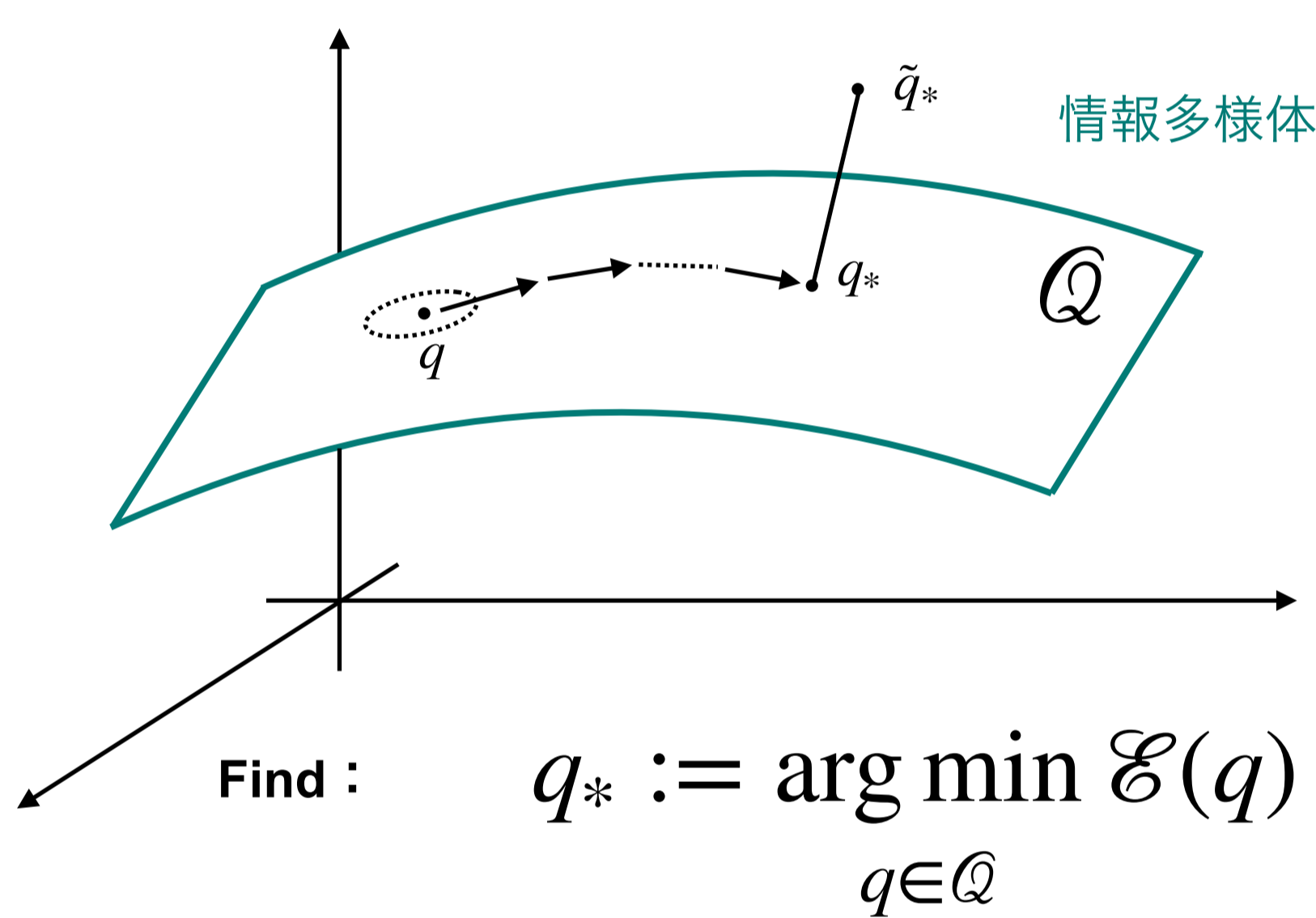
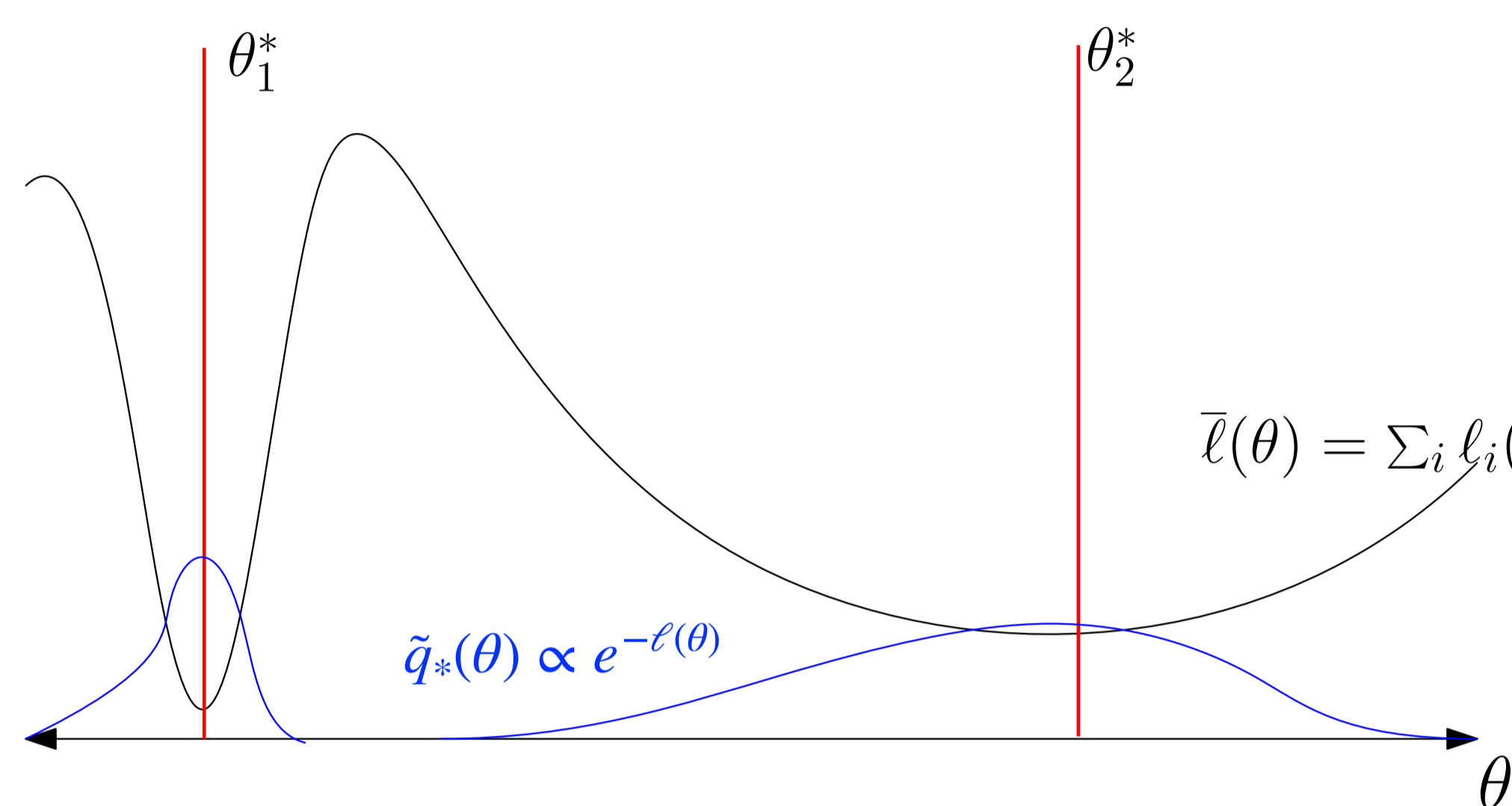
# The Lie Group Bayesian Rule

Eren Mehmet Kiral<sup>1,2</sup>, Thomas Möllenhoff<sup>1</sup>, Emtiyaz Khan<sup>1</sup>

1. RIKEN AIP 2. RIKEN SPDR

AISTATS 2023

教師付き学習とは、パラメーター付の関数が与えられたとき、データに対する損失関数 (loss function) が最小となるパラメーターの値を探す問題である。ベイズ推論 (Bayesian Learning Problem) とは、「値」を「分布」で置き換えて得られる問題である。



Bayesian Update Ruleは、は「良い」分布族を与える情報多様体  $\mathcal{Q}$  の中から、最小分布  $\tilde{q}_*$  を近似する分布  $q_*$  を見つける手法である。荒い解から、情報多様体上のニュートン法の類似の手法を用いて、より良い解を見つけ出す。

従来手法であるBayesian Update Ruleでは、2つの制約があった

- 扱える分布族は、基本的にガウス型分布族だけ
- Bayesian Update Ruleの適用で、パラメーターが  $\mathcal{Q}$  から外れてしまう恐れがある

本研究のLie Group Bayesian Ruleは、分布族のパラメーターとして「Lie群」と呼ばれる代数構造をもつ対象を用いるというアイデアに基づいている。Lie群により幅広い分布族を扱えると同時に、Update RuleもLie群の代数構造を用いることにより、パラメーターはLie群の中に収まる形で定義できる。従来手法であるBayesian Update RuleはLie Group Bayesian Ruleの特殊な場合の線形化と解釈できることも示した。様々な分布族やLie群に対して適用でき、従来とは異なる特性を持つ学習アルゴリズムなども見出すことに成功している。

## Bayesian Learning Problem

$\ell(\theta)$ , a loss function on model parameters  $\theta \in \Theta$ . We solve

$$q_* \in \arg \min_{q \in \mathcal{Q}} \underbrace{\mathbb{E}_q[\ell] - \tau \mathcal{H}(q)}_{=: \mathcal{E}(q)}$$

for some family of distributions  $\mathcal{Q}$  on  $\Theta$ .

- The expectation  $\mathbb{E}_q[\ell] = \int_{\Theta} \ell(\theta) q(\theta) d\theta$  prefers regions of  $\ell$  with low loss.
- The entropy  $\mathcal{H}(q) = - \int_{\Theta} q(\theta) \log q(\theta) d\theta$  prefers a higher spread of  $q$ .
- The temperature  $\tau > 0$ , a balancing term.

## Enter Lie Groups

We solve the two issues above by considering families parametrized by a Lie group  $G$ . Let  $G$  act on  $\Theta$ , then it acts on a probability distribution by the pushforward action (by reparametrizations). For some distribution  $q_0$

$$\mathcal{Q} := \{q_g = \pi(g)q_0 : g \in G\}$$

where  $(\pi(g)q)(\theta) = \left| \frac{d(g \cdot \theta)}{d\theta} \right|^{-1} q(g^{-1} \cdot \theta)$ .

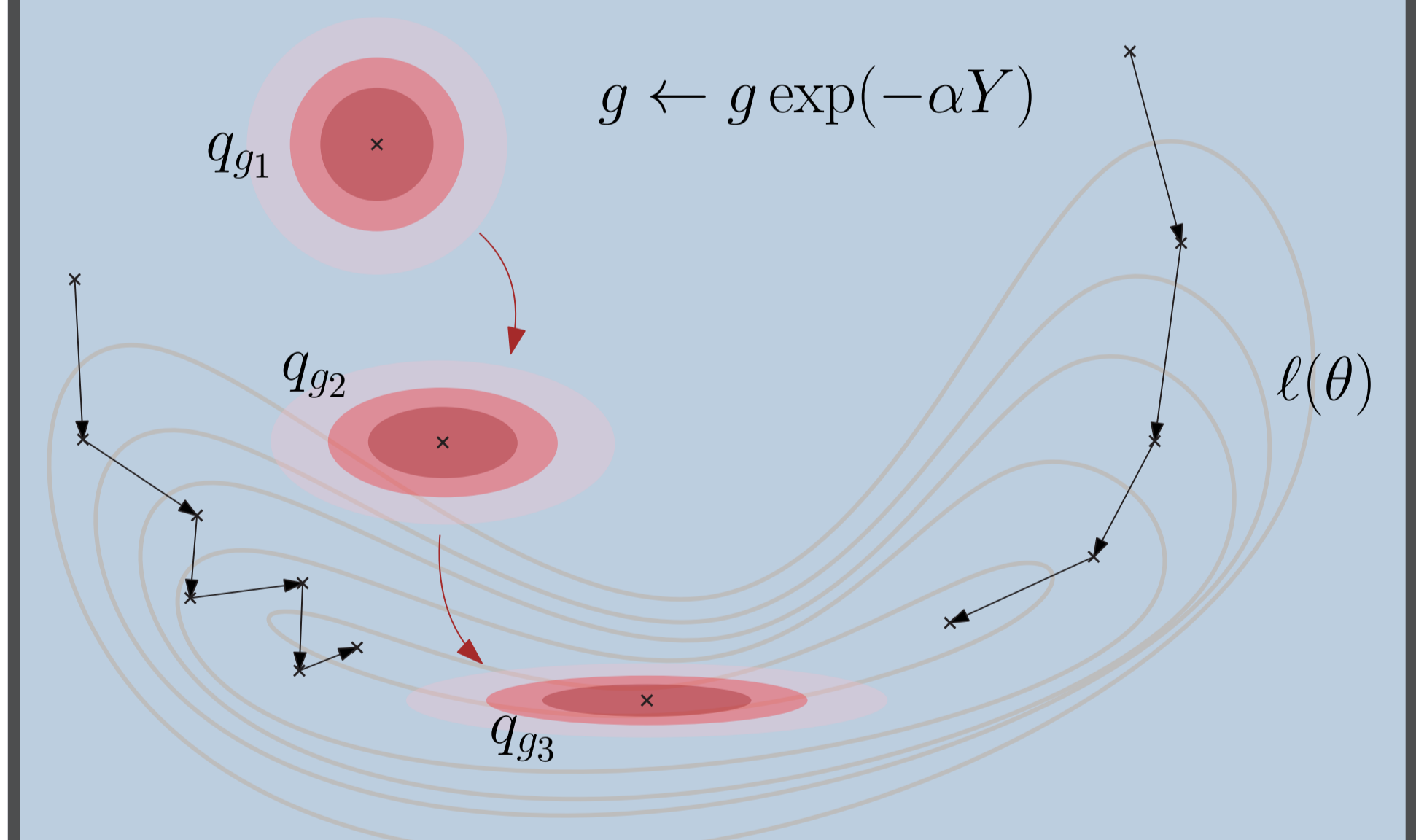
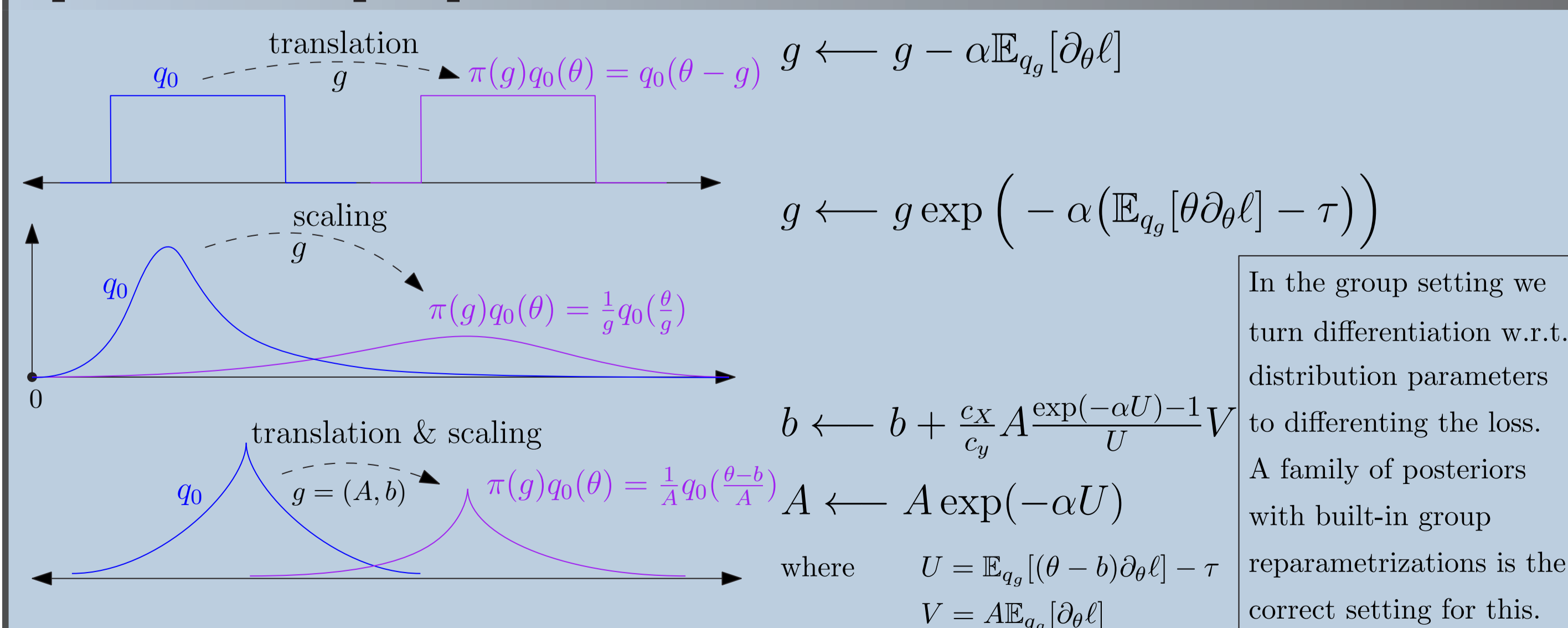


Figure 1: Updating the group element is sufficient to update the distribution.  $Y$  is a vector in the Lie algebra (tangent space of  $G$  at identity) in the direction of steepest ascent of  $\mathcal{E}(q)$  at  $q = q_g$  with respect to the Fisher metric. The update always produces a group element  $g$ , and therefore a distribution in  $\mathcal{Q}$ . And we have flexibility in the choice of  $q_0$ .

## Specific Group Updates



Dataset	Metric	Affine gaussian			Additive [Or22]			[SGD]
		uniform	laplace	laplace	uniform	gaussian	laplace	
CIFAR-10	Acc. (↑)	91.60±0.05	91.53±0.10	91.87±0.04	91.07±0.08	91.28±0.11	91.14±0.12	91.22±0.07
	NLL (↓)	0.300±0.002	0.294±0.004	0.272±0.002	0.365±0.003	0.328±0.008	0.312±0.005	0.354±0.006
	ECE (↓)	0.040±0.001	0.036±0.001	0.029±0.001	0.052±0.001	0.045±0.001	0.039±0.001	0.050±0.001
CIFAR-100	Acc. (↑)	66.08±0.12	66.55±0.10	66.44±0.10	64.29±0.09	64.61±0.20	64.85±0.13	64.19±0.14
	NLL (↓)	1.288±0.007	1.255±0.005	1.247±0.006	1.437±0.004	1.390±0.008	1.359±0.006	1.431±0.007
	ECE (↓)	0.093±0.002	0.079±0.002	0.071±0.001	0.121±0.001	0.107±0.001	0.096±0.001	0.121±0.001
TinyImgNet	Acc. (↑)	51.19±0.12	51.13±0.16	51.36±0.14	49.34±0.14	49.62±0.15	49.73±0.18	49.48±0.10
	NLL (↓)	2.099±0.005	2.098±0.004	2.101±0.009	2.234±0.008	2.204±0.003	2.184±0.007	2.231±0.004
	ECE (↓)	0.076±0.001	0.070±0.002	0.065±0.001	0.107±0.001	0.099±0.002	0.089±0.001	0.106±0.001

## Sparse and localized behavior of the multiplicative rule

Dataset	Metric	additive	multiplicative
MNIST / MLP	Acc. (↑)	98.38±0.02	98.59±0.02
	NLL (↓)	0.083±0.001	0.058±0.001
	ECE (↓)	0.012±0.000	0.006±0.000
CIFAR-10 / MLP	Acc. (↑)	58.85±0.08	59.19±0.07
	NLL (↓)	1.236±0.002	1.160±0.001
	ECE (↓)	0.085±0.001	0.026±0.001

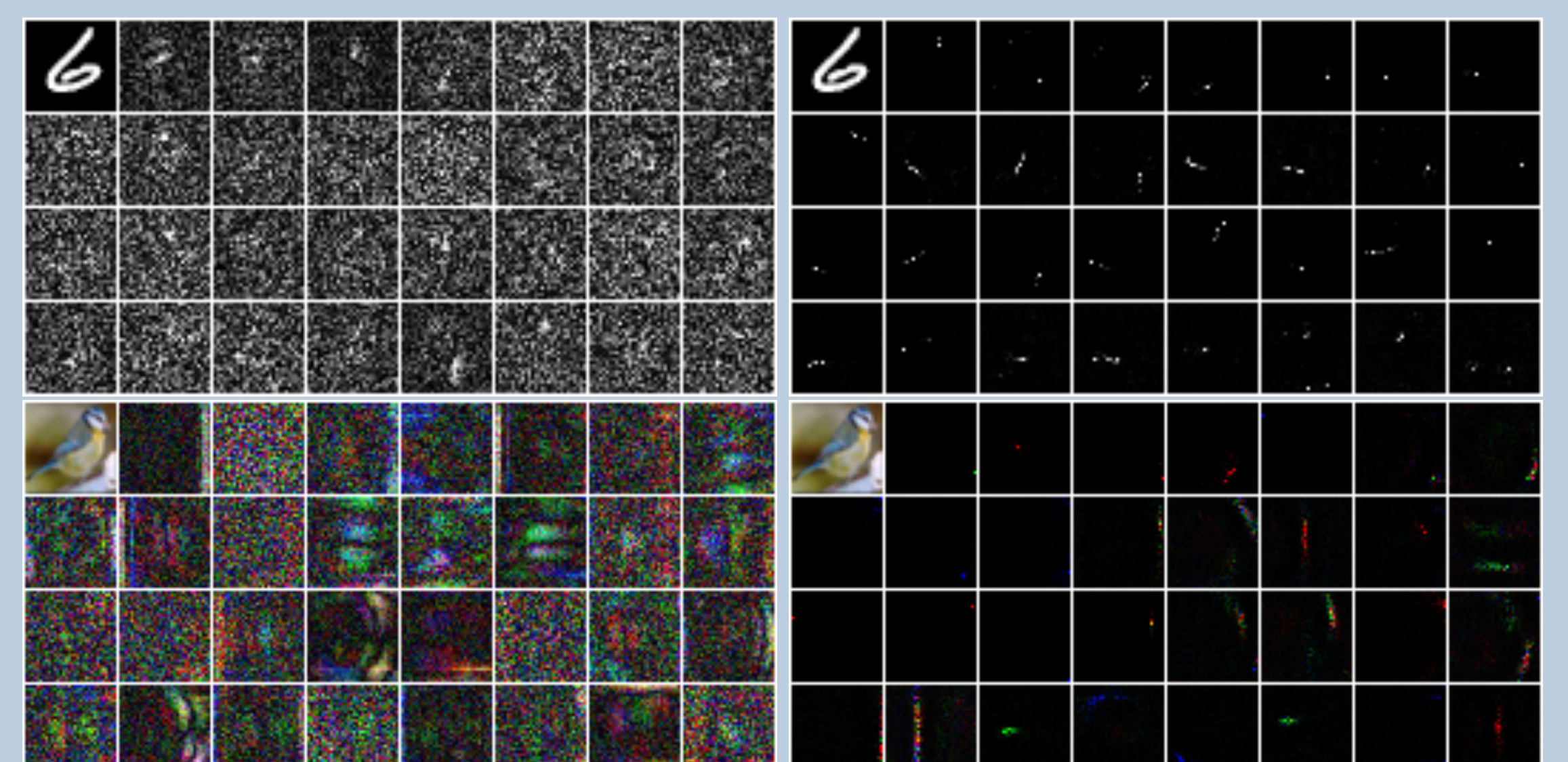


Figure 2: Highest activated first layer filters for the given figure in MNIST (top) or CIFAR10 (bottom); for the additive (left) or multiplicative (right) update rule