**Theorem 1** (Monotonicity of LV–EBM ELBO over variational families)**.** Fix a dataset $\{x_i\}_{i=1}^N$ and a latent-variable energy-based model

$$p_\theta(x,z) \;=\; \frac{1}{Z(\theta)}\exp\big(-E_\theta(x,z)\big),$$

with $Z(\theta) < \infty$ and such that, for each $i$, the conditional density $p_\theta(z \mid x_i)$ exists and is strictly positive on its support.

Let $\mathcal{Q}$ be the set of probability densities $q(z)$ on the latent space with finite differential entropy, and let $\tilde{\mathcal{Q}}$ be the analogous set of joint densities $\tilde{q}(x,z)$ with finite differential entropy. For $q \in \mathcal{Q}$ and $\tilde{q} \in \tilde{\mathcal{Q}}$ define

$$A_i(q) \;:=\; -\mathbb{E}_{q(z)}\big[E_\theta(x_i,z)\big] + H(q), \qquad B(\tilde{q}) \;:=\; -\mathbb{E}_{\tilde{q}(x,z)}\big[E_\theta(x,z)\big] + H(\tilde{q}),$$

and for any collections of admissible families $S_{\text{pos}} = \{S_{\text{pos}}(x_i)\}_{i=1}^N$ with $S_{\text{pos}}(x_i) \subseteq \mathcal{Q}$ and $S_{\text{neg}} \subseteq \tilde{\mathcal{Q}}$ define

$$\mathrm{ELBO}(\theta; S_{\text{pos}}, S_{\text{neg}}) \;:=\; \frac{1}{N}\sum_{i=1}^N \sup_{q \in S_{\text{pos}}(x_i)} A_i(q) \;-\; \sup_{\tilde{q} \in S_{\text{neg}}} B(\tilde{q}).$$

Assume that, for all $\theta$ and all admissible choices of $S_{\text{pos}}, S_{\text{neg}}$, the above suprema are finite. Let

$$S_{\text{pos}}^1 = \big\{S_{\text{pos}}^1(x_i)\big\}_{i=1}^N, \qquad S_{\text{pos}}^2 = \big\{S_{\text{pos}}^2(x_i)\big\}_{i=1}^N$$

be two collections of positive-phase families with $S_{\text{pos}}^1(x_i) \subseteq S_{\text{pos}}^2(x_i) \subseteq \mathcal{Q}$ for all $i$, and let $S_{\text{neg}}^1 \subseteq S_{\text{neg}}^2 \subseteq \tilde{\mathcal{Q}}$ be two admissible negative-phase families. Then, for every $\theta$,

$$\mathrm{ELBO}\big(\theta; S_{\text{pos}}^1, S_{\text{neg}}^1\big) \;\leq\; \mathrm{ELBO}\big(\theta; S_{\text{pos}}^2, S_{\text{neg}}^2\big) \;\leq\; \frac{1}{N}\sum_{i=1}^N \log p_\theta(x_i). \tag{1}$$

Moreover, the first inequality in (1) is strict for a given $\theta$ whenever at least one of the following holds:

(i) For some $i$, every maximizer of $\sup_{q \in S_{\text{pos}}^2(x_i)} A_i(q)$ lies outside $S_{\text{pos}}^1(x_i)$.

(ii) Every maximizer of $\sup_{\tilde{q} \in S_{\text{neg}}^2} B(\tilde{q})$ lies outside $S_{\text{neg}}^1$.

*Proof.* By the Donsker–Varadhan variational formula applied to $f_x(z) := -E_\theta(x,z)$ and $f(x,z) := -E_\theta(x,z)$, we have

$$\log p_\theta(x_i) = \sup_{q \in \mathcal{Q}} A_i(q),$$
$$\log Z(\theta) = \sup_{\tilde{q} \in \tilde{\mathcal{Q}}} B(\tilde{q}),$$

with maximizers $q^\star(\cdot) = p_\theta(\cdot \mid x_i)$ and $\tilde{q}^\star = p_\theta(x,z)$, respectively. Equivalently, using the identities

$$\mathbb{E}_{\tilde{q}}[E_\theta] - H(\tilde{q}) = \mathrm{KL}(\tilde{q}\|p_\theta) - \log Z(\theta), \qquad \mathbb{E}_q[E_\theta(x_i,\cdot)] - H(q) = \mathrm{KL}(q\|p_\theta(\cdot \mid x_i)) - \log p_\theta(x_i),$$

we obtain

$$A_i(q) = \log p_\theta(x_i) - \mathrm{KL}\big(q \parallel p_\theta(\cdot \mid x_i)\big),$$
$$B(\tilde{q}) = \log Z(\theta) - \mathrm{KL}\big(\tilde{q} \parallel p_\theta\big).$$

Therefore, for any $S_{\text{pos}}, S_{\text{neg}}$,

$$\sup_{q \in S_{\text{pos}}(x_i)} A_i(q) = \log p_\theta(x_i) - \inf_{q \in S_{\text{pos}}(x_i)} \mathrm{KL}\big(q \parallel p_\theta(\cdot \mid x_i)\big),$$
$$\sup_{\tilde{q} \in S_{\text{neg}}} B(\tilde{q}) = \log Z(\theta) - \inf_{\tilde{q} \in S_{\text{neg}}} \mathrm{KL}\big(\tilde{q} \parallel p_\theta\big),$$

and hence

$$\mathrm{ELBO}(\theta; S_{\text{pos}}, S_{\text{neg}}) = \frac{1}{N}\sum_{i=1}^N \log p_\theta(x_i) - \log Z(\theta) - \frac{1}{N}\sum_{i=1}^N \inf_{q \in S_{\text{pos}}(x_i)} \mathrm{KL}\big(q \parallel p_\theta(\cdot \mid x_i)\big) + \inf_{\tilde{q} \in S_{\text{neg}}} \mathrm{KL}\big(\tilde{q} \parallel p_\theta\big).$$

1

Now take two nested pairs $(S_{\text{pos}}^1, S_{\text{neg}}^1)$ and $(S_{\text{pos}}^2, S_{\text{neg}}^2)$ as in the statement. For each $i$,

$$\inf_{q \in S_{\text{pos}}^2(x_i)} \text{KL}\big(q \parallel p_\theta(\cdot \mid x_i)\big) \ \leq \ \inf_{q \in S_{\text{pos}}^1(x_i)} \text{KL}\big(q \parallel p_\theta(\cdot \mid x_i)\big),$$

and likewise

$$\inf_{\tilde{q} \in S_{\text{neg}}^2} \text{KL}\big(\tilde{q} \parallel p_\theta\big) \ \leq \ \inf_{\tilde{q} \in S_{\text{neg}}^1} \text{KL}\big(\tilde{q} \parallel p_\theta\big),$$

since infima over supersets cannot be larger. Substituting these inequalities into the expression for $\text{ELBO}(\theta; \cdot, \cdot)$ yields

$$\text{ELBO}\big(\theta; S_{\text{pos}}^1, S_{\text{neg}}^1\big) \ \leq \ \text{ELBO}\big(\theta; S_{\text{pos}}^2, S_{\text{neg}}^2\big)$$

for all $\theta$, with strict inequality whenever at least one of the two infima is strictly smaller over $S_{\text{pos}}^2, S_{\text{neg}}^2$ than over $S_{\text{pos}}^1, S_{\text{neg}}^1$; this is exactly conditions (i)–(ii).

Finally, taking $S_{\text{pos}}^2(x_i) = \mathcal{Q}$ and $S_{\text{neg}}^2 = \tilde{\mathcal{Q}}$, we recover $\sup_{q \in \mathcal{Q}} A_i(q) = \log p_\theta(x_i)$ and $\sup_{\tilde{q} \in \tilde{\mathcal{Q}}} B(\tilde{q}) = \log Z(\theta)$, so

$$\text{ELBO}(\theta; \mathcal{Q}, \tilde{\mathcal{Q}}) = \frac{1}{N} \sum_{i=1}^N \log p_\theta(x_i),$$

which gives the upper bound in (1). $\qquad\square$

**Corollary 1** (LV–EBM gradient-flow families dominate parametric VI)**.** In the setting of Theorem 1, fix $\theta$ and let $S_{\text{pos}}^{\text{VI}}(x_i) \subseteq \mathcal{Q}$ and $S_{\text{neg}}^{\text{VI}} \subseteq \tilde{\mathcal{Q}}$ be any finite-dimensional "variational inference" families (e.g. diagonal Gaussian conditionals and a parametric joint). Define the corresponding VI objective by

$$\text{ELBO}_{\text{VI}}(\theta) \ := \ \text{ELBO}\big(\theta; S_{\text{pos}}^{\text{VI}}, S_{\text{neg}}^{\text{VI}}\big).$$

Let $(\mathcal{T}_t^x)_{t \geq 0}$ and $(\tilde{\mathcal{T}}_t)_{t \geq 0}$ denote the Markov semigroups associated with the conditional and joint Langevin dynamics used in LV–EBMs, assumed to be well-defined and Feller under the regularity assumptions on $E_\theta$ in the main text. For any fixed time horizon $t \geq 0$, define the (infinite-dimensional) gradient-flow families

$$S_{\text{pos}}^{\text{flow}}(x_i) \ := \ \big\{\mathcal{T}_t^{x_i}(q_0) : q_0 \in \mathcal{Q}\big\} \subseteq \mathcal{Q}, \qquad S_{\text{neg}}^{\text{flow}} \ := \ \big\{\tilde{\mathcal{T}}_t(\tilde{q}_0) : \tilde{q}_0 \in \tilde{\mathcal{Q}}\big\} \subseteq \tilde{\mathcal{Q}}.$$

Assume that every VI density can be realized as an admissible initial distribution for the semigroups, and that we include $t = 0$ in the definition of the flow families (so that $\mathcal{T}_0^x$ and $\tilde{\mathcal{T}}_0$ are the identity maps on $\mathcal{Q}$ and $\tilde{\mathcal{Q}}$). Then

$$S_{\text{pos}}^{\text{VI}}(x_i) \ \subseteq \ S_{\text{pos}}^{\text{flow}}(x_i), \qquad S_{\text{neg}}^{\text{VI}} \ \subseteq \ S_{\text{neg}}^{\text{flow}},$$

and, for all $\theta$,

$$\text{ELBO}_{\text{VI}}(\theta) \ \leq \ \text{ELBO}\big(\theta; S_{\text{pos}}^{\text{flow}}, S_{\text{neg}}^{\text{flow}}\big) \ \leq \ \frac{1}{N} \sum_{i=1}^N \log p_\theta(x_i). \tag{2}$$

Moreover, the first inequality in (2) is strict whenever, for the given $\theta$, at least one of the VI families $S_{\text{pos}}^{\text{VI}}(x_i)$ or $S_{\text{neg}}^{\text{VI}}$ is misspecified in the sense that its best-approximation KL divergence to the true posterior $p_\theta(z \mid x_i)$ or joint $p_\theta(x, z)$ is strictly positive.

*Proof.* By definition of the semigroups, $\mathcal{T}_0^x$ and $\tilde{\mathcal{T}}_0$ act as the identity on $\mathcal{Q}$ and $\tilde{\mathcal{Q}}$, so any VI density $q_\phi(z \mid x_i)$ or $q_\psi(x, z)$ can be written as $\mathcal{T}_0^{x_i}(q_0)$ or $\tilde{\mathcal{T}}_0(\tilde{q}_0)$ with $q_0 = q_\phi$ and $\tilde{q}_0 = q_\psi$. Hence $S_{\text{pos}}^{\text{VI}}(x_i) \subseteq S_{\text{pos}}^{\text{flow}}(x_i)$ and $S_{\text{neg}}^{\text{VI}} \subseteq S_{\text{neg}}^{\text{flow}}$. Applying Theorem 1 with $(S_{\text{pos}}^1, S_{\text{neg}}^1) = (S_{\text{pos}}^{\text{VI}}, S_{\text{neg}}^{\text{VI}})$ and $(S_{\text{pos}}^2, S_{\text{neg}}^2) = (S_{\text{pos}}^{\text{flow}}, S_{\text{neg}}^{\text{flow}})$ gives the inequalities in (2). Strictness follows from the strict part of Theorem 1 whenever at least one of the VI families cannot realize the exact posterior or joint. $\qquad\square$