

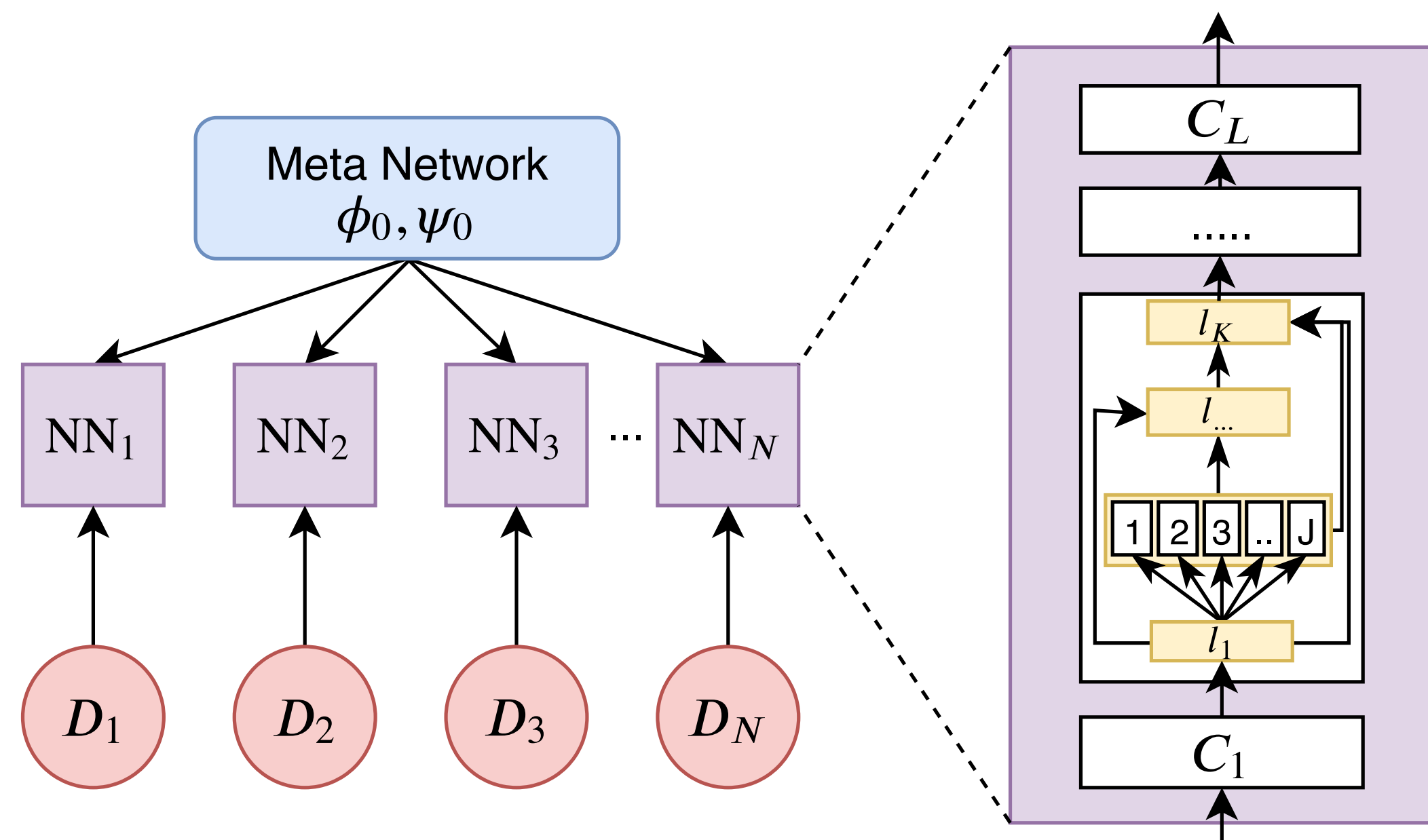
Motivation

While developing techniques to automatically search the space of Neural Network architectures has become a large focus of recent efforts, existing methods can only learn the architecture for a single task at a time. We approach the problem of neural architecture search using Bayesian and Meta-learning techniques to simultaneously learn the architecture and weight distribution for multiple tasks at once.

Main Contribution:

- We propose a **Bayesian inference** view of architecture learning.
- We use this novel view to derive a variational inference method to learn the architecture of a meta-network, which will be shared across multiple tasks using the **optimization embedding** technique to design the parameterization of the posterior.
- Demonstrates a **concrete algorithm** for Meta Architecture Search that can use a prior trained over multiple tasks to find competitive models for new unseen datasets with just quick adaptation.

Bayesian View of Architecture Search



We consider neural network architecture learning as an operation selection problem. Starting with a DenseNet architecture as shown above where the k -th layer of the neural network is defined as

$$x_k = \sum_{i=1}^{k-1} (z_{i,k}^\top \mathcal{A}_i(\theta)) \circ x_i := \sum_{i=1}^{k-1} \sum_{l=1}^L z_{i,k}^l \phi_i^l(x_i; \theta),$$

where $\mathcal{A}_i(\theta) = [\phi_i^l(\cdot; \theta)]_{l=1}^L$ denotes a group of operations from $\mathbb{R}^d \rightarrow \mathbb{R}^p$ depending on the parameters θ and $z_{i,k} \sim \text{Categorical}(\alpha_{i,k})$ with $\alpha_{i,k}^l \geq 0$, $\sum_{l=1}^L \alpha_{i,k}^l = 1$. Assume the probabilistic model as

$$\begin{aligned} \theta &\sim \mathcal{N}(\mu, \sigma^2), \\ z_{i,k} &\sim \text{Categorical}(\alpha_{i,k}), \quad k = 1, \dots, K, \\ y &\sim p(y|x; \theta, z) \propto \exp(-\ell(f(x; \theta, z), y)), \end{aligned}$$

We can estimate the parameters $W(\mu, \sigma, \alpha)$ via maximum log-likelihood estimation

$$\max_W \hat{\mathbb{E}}_{x,y} \left[\log \int p(y|x; \theta, z) p(z; \alpha) p(\theta; \mu, \sigma) dz d\theta \right].$$

However, the MLE is intractable due to the integral. We consider the ELBO, i.e.,

$$\max_W \max_{q(z), q(\theta)} -\hat{\mathbb{E}}_{x,y} \mathbb{E}_{z \sim q(z), \theta \sim q(\theta)} [\ell(f(x; \theta, z), y)] - KL(q||p),$$

whose $p(z) = \prod_{k=1}^K \prod_{i=1}^{k-1} \text{Categorical}(z_{i,k}) = \prod_{k=1}^K \prod_{i=1}^{k-1} \prod_{l=1}^L (\alpha_{i,k}^l)^{z_{i,k}^l}$.

Bayesian Meta-network Architecture Learning

The Bayesian view can be extended to the few-shot meta-learning setting, where we have many tasks, i.e., $\mathcal{D}_t = \{x_i^t, y_i^t\}_{i=1}^n$.

This can be formulated so that the layers and architecture priors i.e., the hyperparameters (μ, σ, α) are shared between all the tasks.

$$\max_W \hat{\mathbb{E}}_{\mathcal{D}_t} \hat{\mathbb{E}}_{(x,y) \sim \mathcal{D}_t} \left[\log \int p(y|x; \theta, z) p(z; \alpha) p(\theta; \mu, \sigma) dz d\theta \right]$$

Using the convexity of log-sum-exp we achieve the ELBO

$$\max_W \hat{\mathbb{E}}_{\mathcal{D}_t} \left[\max_{q(z|\mathcal{D}), q(\theta|\mathcal{D})} \hat{\mathbb{E}}_{(x,y) \sim \mathcal{D}_t} \mathbb{E}_{z \sim q(z|\mathcal{D}), \theta \sim q(\theta|\mathcal{D})} [-\ell(f(x; \theta, z), y)] - KL(q||p) \right]$$

With the variational posterior distributions, $q(z|\mathcal{D})$ and $q(\theta|\mathcal{D})$, introduced into the model, we can directly generate the architecture and its corresponding weights based on the posterior. In a sense, the posterior can be understood as the neural network predictive model.

Variational Inference by Optimization Embedding

The design of the parameterization of the posterior $q(z|\mathcal{D})$ and $q(\theta|\mathcal{D})$ is extremely important, especially in our case where we need to model the dependence w.r.t. the task distributions \mathcal{D} . We assume the $q(\theta|\mathcal{D})$ is Gaussian and the $q(z|\mathcal{D})$ is product of the Categorical distribution which we approximate with Gumbel-Softmax. Therefore, we have

$$q_\psi(\theta|\mathcal{D}) = \mathcal{N}(\psi_\mu, \psi_\sigma), \quad q_\phi(z_{i,k}|\mathcal{D}) = \Gamma(r) \tau^{L-1} \left(\sum_{l=1}^L \frac{\pi_{\mathcal{D}, \phi_{i,k}^l}}{(z_{i,k}^l)^\tau} \right)^{-r} \prod_{l=1}^r \left(\frac{\pi_{\mathcal{D}, \phi_{i,k}^l}}{(z_{i,k}^l)^{\tau+1}} \right),$$

Then, we can sample (θ, z) by following,

$$\begin{aligned} \theta_{\mathcal{D}}(\epsilon, \psi) &= \psi_{\mathcal{D}, \mu} + \epsilon \psi_{\mathcal{D}, \sigma}, \quad \epsilon \sim \mathcal{N}(0, 1), \\ z_{i,k,\mathcal{D}}^l(\xi, \phi) &= \frac{\exp\left(\left(\phi_{\mathcal{D}, i,k}^l + \xi^l\right) / \tau\right)}{\sum_{l=1}^L \exp\left(\left(\phi_{\mathcal{D}, i,k}^l + \xi^l\right) / \tau\right)}, \quad \xi^l \sim \mathcal{G}(0, 1), \quad l \in \{1, \dots, L\}, \end{aligned}$$

with $\pi_{x,\phi,i} = \frac{\exp(\phi_{x,i})}{\sum_{j=1}^p \exp(\phi_{x,j})}$ and $\mathcal{G}(0, 1)$ denoting the Gumbel distribution.

Plugging the formulation into the ELBO, we arrive at the objective

$$\hat{\mathbb{E}}_{\mathcal{D}} \left[\max_{\phi_{\mathcal{D}}, \psi_{\mathcal{D}}} \underbrace{\hat{\mathbb{E}}_{x,y} \mathbb{E}_{\xi, \epsilon} [-\ell(f(x; \theta_{\mathcal{D}}(\epsilon, \psi), z_{\mathcal{D}}(\xi, \phi)), y)]}_{L(\phi_{\mathcal{D}}, \psi_{\mathcal{D}}; W)} - \log \frac{q_\phi(z|\mathcal{D})}{p(z; \alpha)} - \log \frac{q_\psi(\theta|\mathcal{D})}{p(\theta; \mu, \sigma)} \right].$$

We follow the parametrized CVB derivation for embedding the optimization procedure for (ϕ, ψ) deriving the explicit form of the parameters $\phi_{\mathcal{D}}$ and $\psi_{\mathcal{D}}$.

Denoting the $\hat{g}_{\phi_{\mathcal{D}}, \psi_{\mathcal{D}}}(\mathcal{D}, W) = \frac{\partial \hat{L}}{\partial (\phi_{\mathcal{D}}, \psi_{\mathcal{D}})}$ where \hat{L} is the stochastic approximation for $L(\phi_{\mathcal{D}}, \psi_{\mathcal{D}}; W)$, we have

$$[\phi_{\mathcal{D}}^t, \psi_{\mathcal{D}}^t] = \eta_t \hat{g}_{\phi_{\mathcal{D}}, \psi_{\mathcal{D}}}(\mathcal{D}, W) + [\phi_{\mathcal{D}}^{t-1}, \psi_{\mathcal{D}}^{t-1}],$$

We can initialize $(\phi^0, \psi^0) = W$ which is shared across all the tasks. After T optimization steps, we obtain $(\phi_{\mathcal{D}}^T, \psi_{\mathcal{D}}^T)$, which leads to $(\theta_{\mathcal{D}}^T(\xi, \psi_{\mathcal{D}}^T), z_{\mathcal{D}}(\xi, \phi_{\mathcal{D}}^T))$. In other words, we derive the concrete parameterization of $q(\theta|\mathcal{D})$ and $q(z|\mathcal{D})$ automatically by unfolding the optimization steps. Plugging the ultimate parameterization into $L(\phi_{\mathcal{D}}, \psi_{\mathcal{D}}; W)$, we have

$$\max_{W,V} \hat{\mathbb{E}}_{\mathcal{D}} \hat{\mathbb{E}}_{x,y} \mathbb{E}_{\xi, \epsilon} \left[-\ell(f(x; \theta_{\mathcal{D}}^T(\epsilon, \psi), z_{\mathcal{D}}^T(\xi, \phi)), y) - \log \frac{q_{\phi_{\mathcal{D}}^T}(z|\mathcal{D})}{p(z; \alpha)} - \log \frac{q_{\psi_{\mathcal{D}}^T}(\theta|\mathcal{D})}{p(\theta; \mu, \sigma)} \right],$$

which can be optimized by stochastic gradient ascent for learning W .

Bayesian meta-Architecture SEarch (BASE) Algorithm

- 1: Initialize meta-network parameters W_0 .
- 2: **for** $e = 1, \dots, E$ **do**
- 3: Sample C tasks $\{\mathcal{D}_c\}_{c=1}^C \sim \mathcal{D}$.
- 4: **for** \mathcal{D}_c in \mathcal{D} **do**
- 5: Sample $\{x_t, y_t\}_{t=1}^T \sim \mathcal{D}_c$.
- 6: Let $\phi_c^0, \psi_c^0 = W_{e-1}$.
- 7: **for** $t = 1, \dots, T$ **do**
- 8: Sample $\xi \sim \mathcal{G}(0, 1)$.
- 9: Update $[\phi_c^t, \psi_c^t] = [\phi_c^{t-1}, \psi_c^{t-1}] - \eta \nabla_{\phi_c^{t-1}, \psi_c^{t-1}} \hat{L}(f(x_t; \phi_c^{t-1}, \psi_c^{t-1}, \xi), y_t)$.
- 10: Update $W_e = W_{e-1} + \lambda \frac{1}{C} \sum_{c=1}^C ([\phi_c^T, \psi_c^T] - W_{e-1})$.

Classification Accuracy on CIFAR10

Architecture	Top-1 Test Error	Params (M)	Search Time (Gpu Days)
DenseNet-BC	3.46	25.6	-
NASNet-A + cutout	2.65	3.3	1800
AmoebaNet-A + cutout	3.34 ± 0.06	3.2	3150
AmoebaNet-B + cutout	2.55 ± 0.05	2.8	3150
Hierarchical Evo	3.75 ± 0.12	15.7	300
DARTS (1st order)	3.00 ± 0.14	3.3	1.5
DARTS (2nd order)	2.76 ± 0.09	3.3	4
SNAS (single-level)	2.85 ± 0.02	2.8	1.5
ENAS + cutout	2.89	4.6	0.5
PNAS	3.41 ± 0.09	3.2	225
SMASH	4.03	16	1.5
BASE(Multi-task Prior)	3.18	3.22	8
BASE(Imagenet32)	3.00	3.29	0.04 Adap / 8 Meta
BASE(CIFAR10)	2.83	3.07	0.05 Adap / 8 Meta

Classification Accuracy on Imagenet

Architecture	Top-1 Err	Top-5 Err	MACs (M)	Search Time (GPU Days)
NASNet-A	26.0	8.4	564	1800
NASNet-B	27.2	8.7	488	1800
NASNet-C	27.5	9.0	558	1800
AmoebaNet-A	25.5	8.0	555	3150
AmoebaNet-B	26.0	8.5	555	3150
AmoebaNet-C	24.3	7.6	570	3150
PNAS	25.8	8.1	588	225
DARTS	26.9	9.0	595	4
SNAS	27.3	9.2	522	1.5
BASE (Multi-task Prior)	26.12	8.52	544	0 Adap / 8 Meta
BASE (Imagenet)	25.71	8.08	559	0.04 Adap / 8 Meta