



Meta-learning Adaptive Deep Kernel Gaussian Processes for Molecular Property Prediction

Wenlin Chen, Austin Tripp, José Miguel Hernández-Lobato

Published at ICLR 23: <u>https://openreview.net/forum?id=KXRSh0sdVTP</u>

Goal: find novel molecules with desirable biochemical/physicochemical properties.



Goal: find novel molecules with desirable biochemical/physicochemical properties.



Challenge:

- Evaluating molecular properties is slow and expensive: small datasets are ubiquitous.
- Chemical space is huge and complex: exhaustive search is prohibitive.

Goal: find novel molecules with desirable biochemical/physicochemical properties.



Challenge:

- Evaluating molecular properties is slow and expensive: small datasets are ubiquitous.
- Chemical space is huge and complex: exhaustive search is prohibitive.

Gaussian processes (GPs):

- ✓ GPs are well-calibrated models with generally reliable uncertainty on small datasets.
- ✓ GPs could be used as surrogate models in Bayesian optimization to guide molecule search.

Goal: find novel molecules with desirable biochemical/physicochemical properties.



Challenge:

- Evaluating molecular properties is slow and expensive: small datasets are ubiquitous.
- Chemical space is huge and complex: exhaustive search is prohibitive.

Gaussian processes (GPs):

- ✓ GPs are well-calibrated models with generally reliable uncertainty on small datasets.
- ✓ GPs could be used as surrogate models in Bayesian optimization to guide molecule search.
- ★ Hand-designing kernels for structured data like molecules is challenging.

Goal: find novel molecules with desirable biochemical/physicochemical properties.



Challenge:

- Evaluating molecular properties is slow and expensive: small datasets are ubiquitous.
- Chemical space is huge and complex: exhaustive search is prohibitive.

Gaussian processes (GPs):

- ✓ GPs are well-calibrated models with generally reliable uncertainty on small datasets.
- ✓ GPs could be used as surrogate models in Bayesian optimization to guide molecule search.
- ★ Hand-designing kernels for structured data like molecules is challenging.

We want better GP models for molecules!

Model: Deep Kernel Gaussian Processes

Why not just learn features for molecules using a deep neural network (DNN)?

Model: Deep Kernel Gaussian Processes

Why not just learn features for molecules using a deep neural network (DNN)?

Deep kernel GPs operate on features learned by a DNN.



Representation learning (DNN) + Uncertainty (GP)

Model: Deep Kernel Gaussian Processes

Why not just learn features for molecules using a deep neural network (DNN)?

Deep kernel GPs operate on features learned by a DNN.



Representation learning (DNN) + Uncertainty (GP)

learnable parameters $\psi = (\theta, \phi)$ (e.g., RBF) extractor

1. Deep Kernel Learning (DKL, Wilson et al., 2016)

- All parameters are learned by minimizing the negative log marginal likelihood (NLML) on a single dataset. $\psi^* = \arg \min \text{NLML}(\psi, S_{\tau})$

$$\psi^* = \operatorname*{arg\,min}_{\psi} \operatorname{NLML}(\psi, \mathcal{S}_{\mathcal{T}})$$

- Pure single-task learning (a separate deep kernel GP is trained for each task).

Wilson, Andrew Gordon, et al. "Deep kernel learning." Artificial intelligence and statistics. PMLR, 2016.

Ober, Sebastian W., Carl E. Rasmussen, and Mark van der Wilk. "The promises and pitfalls of deep kernel learning." Uncertainty in Artificial Intelligence. PMLR, 2021.

1. Deep Kernel Learning (DKL, Wilson et al., 2016)

- All parameters are learned by minimizing the negative log marginal likelihood (NLML) on a single dataset. $\psi^* = \arg \min \text{NLML}(\psi, S_{\tau})$

$$oldsymbol{\psi}^* = rg\min_{oldsymbol{\psi}} \operatorname{NLML}\left(oldsymbol{\psi}, \mathcal{S}_{\mathcal{T}}
ight)$$

- Pure single-task learning (a separate deep kernel GP is trained for each task).

Severe overfitting (Ober et al., 2021) on **small datasets** despite the use of type-II ML.



DKL makes all output values strongly correlated!

Wilson, Andrew Gordon, et al. "Deep kernel learning." Artificial intelligence and statistics. PMLR, 2016.

Ober, Sebastian W., Carl E. Rasmussen, and Mark van der Wilk. "The promises and pitfalls of deep kernel learning." Uncertainty in Artificial Intelligence. PMLR, 2021.

2. Deep Kernel Transfer (DKT, Patacchiola et al., 2020)

- All parameters are learned by minimizing the expected NLML over a distribution of datasets.

$$oldsymbol{\psi}^* = rgmin_{oldsymbol{\psi}} \mathbb{E}_{p(\mathcal{T})}[\mathrm{NLML}(oldsymbol{\psi},\mathcal{T})]$$

- Pure meta-learning (all parameters are shared across tasks).

2. Deep Kernel Transfer (DKT, Patacchiola et al., 2020)

- All parameters are learned by minimizing the expected NLML over a distribution of datasets.

$$oldsymbol{\psi}^* = rgmin_{oldsymbol{\psi}} \mathbb{E}_{p(\mathcal{T})}[ext{NLML}(oldsymbol{\psi},\mathcal{T})]$$

- Pure meta-learning (all parameters are shared across tasks).

× Underfitting due to model mis-specification.



It's unrealistic to assume all datasets are drawn from an identical GP with the same noise, signal variance, characteristic lengthscales!

- **ADKF-IFT** interpolates between DKL and DKT:

- **ADKF-IFT** interpolates between DKL and DKT:

- Partition the deep kernel parameters $\boldsymbol{\psi} = [\boldsymbol{\theta}, \boldsymbol{\phi}]$ into two disjoint sets $\boldsymbol{\psi}_{meta} = \boldsymbol{\phi}$ and $\boldsymbol{\psi}_{adapt} = \boldsymbol{\theta}$.

- **ADKF-IFT** interpolates between DKL and DKT:

- Partition the deep kernel parameters $\boldsymbol{\psi} = [\boldsymbol{\theta}, \boldsymbol{\phi}]$ into two disjoint sets $\boldsymbol{\psi}_{meta} = \boldsymbol{\phi}$ and $\boldsymbol{\psi}_{adapt} = \boldsymbol{\theta}$.

- Adapt base kernel parameters $\psi_{adapt} = \theta$ to each task's training set S_T by minimizing the NLML \mathcal{L}_T train loss.

$$\psi^*_{\text{adapt}}(\psi_{\text{meta}}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_T(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}}). \quad \longleftarrow \quad \begin{array}{l} \text{best response function for a} \\ \text{given task } \mathcal{T} \text{ and } \psi_{\text{meta}} \end{array}$$

- **ADKF-IFT** interpolates between DKL and DKT:

- Partition the deep kernel parameters $\boldsymbol{\psi} = [\boldsymbol{\theta}, \boldsymbol{\phi}]$ into two disjoint sets $\boldsymbol{\psi}_{meta} = \boldsymbol{\phi}$ and $\boldsymbol{\psi}_{adapt} = \boldsymbol{\theta}$.

- Adapt base kernel parameters $\psi_{adapt} = \theta$ to each task's training set S_T by minimizing the NLML \mathcal{L}_T train loss.
- Meta-learn feature extractor parameters $\psi_{meta} = \phi$ to optimize the model's average performance on the test sets Q_T of many tasks (after $\psi_{adapt} = \theta$ has been separately adapted to each of these tasks).

$$\psi_{\text{meta}}^{*} = \underset{\psi_{\text{meta}}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{\text{meta}},\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$

s.t. $\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{\text{meta}},\psi_{\text{adapt}},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- **ADKF-IFT** interpolates between DKL and DKT:

- Partition the deep kernel parameters $\boldsymbol{\psi} = [\boldsymbol{\theta}, \boldsymbol{\phi}]$ into two disjoint sets $\boldsymbol{\psi}_{meta} = \boldsymbol{\phi}$ and $\boldsymbol{\psi}_{adapt} = \boldsymbol{\theta}$.

- Adapt base kernel parameters $\psi_{adapt} = \theta$ to each task's training set S_T by minimizing the NLML \mathcal{L}_T train loss.
- Meta-learn feature extractor parameters $\psi_{meta} = \phi$ to optimize the model's average performance on the test sets Q_T of many tasks (after $\psi_{adapt} = \theta$ has been separately adapted to each of these tasks).
- The validation loss \mathcal{L}_V is the negative log joint predictive posterior on the test set \mathcal{Q}_T given the training set \mathcal{S}_T .

$$\psi_{\text{meta}}^{*} = \underset{\psi_{\text{meta}}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{\text{meta}},\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$

s.t. $\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{\text{meta}},\psi_{\text{adapt}},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- **ADKF-IFT** interpolates between DKL and DKT:

- Partition the deep kernel parameters $\boldsymbol{\psi} = [\boldsymbol{\theta}, \boldsymbol{\phi}]$ into two disjoint sets $\boldsymbol{\psi}_{meta} = \boldsymbol{\phi}$ and $\boldsymbol{\psi}_{adapt} = \boldsymbol{\theta}$.

- Adapt base kernel parameters $\psi_{adapt} = \theta$ to each task's training set S_T by minimizing the NLML \mathcal{L}_T train loss.
- Meta-learn feature extractor parameters $\psi_{meta} = \phi$ to optimize the model's average performance on the test sets Q_T of many tasks (after $\psi_{adapt} = \theta$ has been separately adapted to each of these tasks).
- The validation loss \mathcal{L}_V is the negative log joint predictive posterior on the test set \mathcal{Q}_T given the training set \mathcal{S}_T .
- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{\text{meta}}^{*} = \underset{\psi_{\text{meta}}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{\text{meta}},\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$

s.t. $\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{\text{meta}},\psi_{\text{adapt}},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Interpretation: DNN meta-learns generally useful features across tasks, such that a task-specific GP operates on top of such features achieves the highest predictive performance on average.

Contrast DKL, DKT and ADKF-IFT



Contrast DKL, DKT and ADKF-IFT



Justification: two related tasks are more likely to have different noise levels, signal variances, or characteristic lengthscales than to require substantially different feature representations.

Contrast DKL, DKT and ADKF-IFT



- ADKF-IFT reduces overfitting:

✓ It regularizes the feature extractor using meta-learning.

✓ It learns feature extractor and base kernel parameters on different subsets (train/test) of a dataset.

- ADKF-IFT reduces underfitting:

✓ It adapts base kernel parameters separately to each task.

Justification: two related tasks are more likely to have different noise levels, signal variances, or characteristic lengthscales than to require substantially different feature representations.

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{\text{meta}}^{*} = \underset{\psi_{\text{meta}}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{\text{meta}},\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$

s.t. $\psi_{\text{adapt}}^{*}(\psi_{\text{meta}},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{\text{meta}},\psi_{\text{adapt}},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{\text{meta}}^{*} = \underset{\psi_{\text{meta}}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{\text{meta}}, \psi_{\text{adapt}}^{*}(\psi_{\text{meta}}, \mathcal{S}_{\mathcal{T}}), \mathcal{T})],$$

s.t. $\psi_{\text{adapt}}^{*}(\psi_{\text{meta}}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{\text{adapt}}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Inner optimization: the gradient of the train loss \mathcal{L}_T can be calculated using auto-diff.

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta}, \psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}), \mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{meta}, \psi_{adapt}, \mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Inner optimization: the gradient of the train loss \mathcal{L}_T can be calculated using auto-diff.

- Outer optimization: how to calculate the gradient for the validation loss \mathcal{L}_V ?

Hypergradient:
$$\frac{d \mathcal{L}_V}{d \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}^*} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}}, \quad \text{(by chain rule)}$$
$$\frac{\text{easy}}{\partial \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}},$$

Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta}, \psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}), \mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg min}} \mathcal{L}_{T}(\psi_{meta}, \psi_{adapt}, \mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Inner optimization: the gradient of the train loss \mathcal{L}_T can be calculated using auto-diff.

- Outer optimization: how to calculate the gradient for the validation loss \mathcal{L}_V ?

Hypergradient:
$$\frac{d \mathcal{L}_V}{d \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}^*} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}}, \quad \text{(by chain rule)}$$
$$\frac{\text{easy}}{\partial \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}^*} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}},$$

The best response function ψ^*_{adapt} is defined by an **argmin function!** How to differentiate it? Auto-diff requires tracking the gradients through many iterations of the inner optimization (**intractable**)!

Solve the Bilevel Optimization Problem by Implicit Function Theorem

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta}, \psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}), \mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{meta}, \psi_{adapt}, \mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Outer optimization: how to calculate the gradient for the validation loss \mathcal{L}_V ?

Hypergradient:
$$\frac{d\mathcal{L}_V}{d\psi_{\text{meta}}} = \frac{\partial\mathcal{L}_V}{\partial\psi_{\text{meta}}} + \frac{\partial\mathcal{L}_V}{\partial\psi_{\text{adapt}}^*} \frac{\partial\psi_{\text{adapt}}^*}{\partial\psi_{\text{meta}}}, \quad \text{(by chain rule)}$$
$$\frac{\text{easy}}{\partial\psi_{\text{meta}}} = \frac{\partial\mathcal{L}_V}{\partial\psi_{\text{adapt}}} \frac{\partial\psi_{\text{adapt}}^*}{\partial\psi_{\text{meta}}},$$

- Since ψ^*_{adapt} is a critical point of the train loss \mathcal{L}_T , we can apply the Implicit Function Theorem (IFT)!

• 1• • 0

Solve the Bilevel Optimization Problem by Implicit Function Theorem

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta},\psi_{adapt}^{*}(\psi_{meta},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{meta},\psi_{adapt},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

9

- Outer optimization: how to calculate the gradient for the validation loss \mathcal{L}_V ?

Hypergradient:
$$\frac{d \mathcal{L}_V}{d \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}^*} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}}, \quad \text{(by chain rule)}$$
$$\frac{\text{easy}}{\text{easy}} = \frac{\text{easy}}{\text{hard}},$$

- Since ψ^*_{adapt} is a critical point of the train loss \mathcal{L}_T , we can apply the Implicit Function Theorem (IFT)!

$$\begin{aligned} \text{IFT:} \quad \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}} \Big|_{\psi_{\text{meta}}'} &= -\left(\frac{\partial^2 \mathcal{L}_T(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}'})}{\partial \psi_{\text{adapt}} \partial \psi_{\text{adapt}}^T}\right)^{-1} \frac{\partial^2 \mathcal{L}_T(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}'})}{\partial \psi_{\text{adapt}} \partial \psi_{\text{meta}}^T} \Big|_{\psi_{\text{meta}}', \psi_{\text{adapt}}'}, \\ \text{(inverse Hessian)} \qquad (\text{mixed partial derivatives}) \end{aligned}$$

Exact and Efficient Gradient Computation



Exact and Efficient Gradient Computation

- Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta},\psi_{adapt}^{*}(\psi_{meta},\mathcal{S}_{\mathcal{T}}),\mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta},\mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{meta},\psi_{adapt},\mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Inner optimization: the gradient of the train loss \mathcal{L}_T can be calculated using auto-diff.

✓ No backpropagate through the feature extractor is required!

✓ Use L-BFGS for base kernel parameter $\psi_{adapt} = \theta$ optimization (fast and efficient).

Exact and Efficient Gradient Computation

Meta-training: ADKF-IFT can be formalized as a bi-level optimization problem.

$$\psi_{meta}^{*} = \underset{\psi_{meta}}{\operatorname{arg\,min}} \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_{V}(\psi_{meta}, \psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}), \mathcal{T})],$$
s.t. $\psi_{adapt}^{*}(\psi_{meta}, \mathcal{S}_{\mathcal{T}}) = \underset{\psi_{adapt}}{\operatorname{arg\,min}} \mathcal{L}_{T}(\psi_{meta}, \psi_{adapt}, \mathcal{S}_{\mathcal{T}}).$ \longleftarrow best response function for a given task \mathcal{T} and ψ_{meta}

- Inner optimization: the gradient of the train loss \mathcal{L}_T can be calculated using auto-diff.

 \checkmark No backpropagate through the feature extractor is required!

✓ Use L-BFGS for base kernel parameter $\psi_{adapt} = \theta$ optimization (fast and efficient).

- Outer optimization: the hypergradient of the validation loss \mathcal{L}_V can be obtained using IFT.

IFT:
$$\frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}}\Big|_{\psi_{\text{meta}}'} = -\left[\frac{\partial^2 \mathcal{L}_T(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}'})}{\partial \psi_{\text{adapt}} \partial \psi_{\text{adapt}}^T}\right]^{-1} \frac{\partial^2 \mathcal{L}_T(\psi_{\text{meta}}, \psi_{\text{adapt}}, \mathcal{S}_{\mathcal{T}'})}{\partial \psi_{\text{adapt}} \partial \psi_{\text{meta}}^T}\Big|_{\psi_{\text{meta}}', \psi_{\text{adapt}}', \psi_{\text{adapt}}'}, \psi_{\text{adapt}}', \psi_{\text{adapt$$

✓ Common GP based kernels (e.g., RBF) contains only a handful of parameters ψ_{adapt} = θ.
 ✓ The inverse Hessian in IFT can be computed exactly without any approximation!

General Framework vs. Specific Instantiations

ADKF-IFT can be formalized as a **bi-level optimization** problem:

$$egin{aligned} &\psi^*_{ ext{meta}} &= rgmin_{m{\psi}_{ ext{meta}}} & \mathbb{E}_{p(\mathcal{T})}[\mathcal{L}_V(\psi_{ ext{meta}},\psi^*_{ ext{adapt}}(\psi_{ ext{meta}},\mathcal{S}_{\mathcal{T}}),\mathcal{T})], \ &\psi^*_{ ext{adapt}}(\psi_{ ext{meta}},\mathcal{S}_{\mathcal{T}}) &= rgmin_{m{\psi}_{ ext{adapt}}} & \mathcal{L}_T(\psi_{ ext{meta}},\psi_{ ext{adapt}},\mathcal{S}_{\mathcal{T}}). \end{aligned}$$

 $\checkmark \psi_{adapt}, \psi_{meta}, \mathcal{L}_T, \mathcal{L}_V$ could be anything, which makes ADKF-IFT a general framework.

✓ Any particular choice of ψ_{adapt} , ψ_{meta} , \mathcal{L}_T , \mathcal{L}_V is an **instantiation** of the general framework.

✓ DKL and DKT are special instantiations (extreme cases) of this general framework!

The General Framework Unifies Previous Methods (DKL and DKT)



(NLML: negative log marginal likelihood)

Experiment 1: Few-shot Molecular Property Prediction on FS-Mol

- FS-Mol (Stanley et al., 2021): 4,938 training tasks, 40 validation tasks, 157 test tasks; 233,786 unique compounds.



✓ The improvements of ADKF-IFT over other methods are statistically significant!

Experiment 1: Ablation Study and Analysis on FS-Mol



$\checkmark \mathbf{DKT} \approx \mathbf{DKT} + \leq \mathbf{ADKF} < \mathbf{ADKF} +$

- **DKT**+ is like **DKT** but tuning the base kernel parameters θ for each task during meta-testing.

- ADKF is like ADKF+ but ignoring the gradient through the best response function ψ^*_{adapt} .

$$\frac{d \mathcal{L}_V}{d \psi_{\text{meta}}} = \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{meta}}} + \frac{\partial \mathcal{L}_V}{\partial \psi_{\text{adapt}}^*} \frac{\partial \psi_{\text{adapt}}^*}{\partial \psi_{\text{meta}}},$$

Experiment 1: Ablation Study and Analysis on FS-Mol



- Blue histogram: the distribution of the base kernel parameters $\boldsymbol{\theta}$ across different tasks learned by ADKF-IFT.

- Dotted vertical line: the base kernel parameters $\boldsymbol{\theta}$ shared across all tasks learned by DKT.

 \checkmark The base kernel parameters θ do vary across tasks!

✓ ADKF-IFT achieves better signal-to-noise ratio!

Experiment 2: OOD Molecular Property Prediction and Optimization

- Bayesian optimization (BO)



- **Surrogate model:** GP operating on top of the features extracted by DNNs meta-trained on FS-Mol by different methods.

- **Evaluation:** four OOD molecular design tasks outside of FS-Mol.

- Test predictive negative log likelihood (NLL)

Feature	Out-of-domain molecular design task							
representation	Molecular docking	Antibiotic discovery	Antiviral drug design	Material design				
Fingerprint	1.138 ± 0.014	1.669 ± 0.075	4.601 ± 0.086	1.091 ± 0.011				
Ρ̈́Α̈́R	1.270 ± 0.019	2.185 ± 0.115	4.840 ± 0.086	1.283 ± 0.017				
MAT	1.528 ± 0.028	2.390 ± 0.104	4.797 ± 0.088	2.198 ± 0.063				
GNN-MT	1.994 ± 0.050	3.692 ± 0.225	6.399 ± 0.181	7.254 ± 0.217				
CNP	1.493 ± 0.028	2.537 ± 0.162	5.005 ± 0.086	1.741 ± 0.043				
ProtoNet	1.147 ± 0.013	1.615 ± 0.094	5.060 ± 0.086	1.032 ± 0.009				
DKT	1.167 ± 0.012	1.602 ± 0.073	4.975 ± 0.092	1.026 ± 0.009				
ADKF-IFT	$\boldsymbol{1.137 \pm 0.011}$	1.496 ± 0.043	4.781 ± 0.087	$\boldsymbol{0.996 \pm 0.007}$				

✓ ADKF-IFT enables fastest discovery of top performing molecules!

✓ ADKF-IFT achieves competitive test predictive performance!

Summary

Our proposed Adaptive Deep Kernel Fitting with Implicit Function Theorem (ADKF-IFT) approach:

- ✓ meta-learns feature representations that facilitate the adaptation of task-specific GP models;
- ✓ generalizes DKL and DKT for training deep kernel GPs using a bilevel optimization framework;
- \checkmark efficiently solve the bilevel optimization problem by implicit function theorem;
- ✓ produces state-of-the-art results on few-shot molecular property prediction benchmarks;
- ✓ achieves great performance on OOD molecular property prediction and optimization tasks;

✓ produces well-calibrated models for fully-automated high-throughput experimentation that could accelerate drug discovery and material design.

Thank you!

Published at ICLR 23: <u>https://openreview.net/forum?id=KXRSh0sdVTP</u>

Limitations and Future Work Directions

1. Use ARD for the lengthscale parameter in the base kernel for automatic feature selection.

- 2. Adapt the feature extractor to each task by allowing small deviations from a meta-learned prior.
- 3. Adopt a more principled approximate inference method for GP classification.

4. Inject domain expertise from drug discovery into the base kernel with hand-curated features and kernel combination.

5. Consider other application domains such as few-shot image classification.

Appendix 1: Mean ranks of Compared Methods on FS-Mol

Table 4: Mean ranks of all compared methods in terms of their performance on all FS-Mol test tasks.

	Support set size						(b) Reg	ression	111 tas	KS).	
Method	16	32	64	128	256			Su	pport set s	ize	
GNN-ST kNN	$11.29 \\ 10.89$	$11.53 \\ 10.48$	$11.75 \\ 10.33$	$\frac{11.85}{10.15}$	$12.19 \\ 9.37$	Method	16	32	64	128	256
MAT RF	$ \begin{array}{r} 10.43 \\ 8.15 \end{array} $	$10.44 \\ 7.89$	$10.19 \\ 7.06$	$9.69 \\ 6.25$	$9.70 \\ 4.47$	MAT GNN-MT	7.60	7.45 6.40	$7.26 \\ 6.15$	$7.06 \\ 5.95$	7.19 5.58
PAR GNN-MT	$7.70 \\ 7.33$	$7.98 \\ 7.18$	$8.30 \\ 7.08$	$8.83 \\ 6.59$	$\begin{array}{c} 10.81 \\ 6.53 \end{array}$	RF	5.00	4.47	4.16	3.72	3.56
DKL GP-ST	$\begin{array}{c} 7.28 \\ 6.71 \end{array}$	$7.49 \\ 6.57$	$7.98 \\ 6.28 \\ 100 \\ 10$	$\begin{array}{c} 8.42 \\ 6.18 \end{array}$	$\begin{array}{c} 8.21 \\ 5.14 \end{array}$	GP-ST	4.42	4.14	3.87	3.37	3.07
GNN-MAML CNP	$6.36 \\ 5.00 \\ 4.00$	$6.92 \\ 5.81 \\ 2.40$	$7.42 \\ 6.36 \\ 2.11$	$7.89 \\ 6.91 \\ 2.00 \\ 0.01 \\ 0.00 \\ $	$8.90 \\ 7.78 \\ 2.95$	DKT	3.88 2.12	$4.45 \\ 2.08$	$4.95 \\ 2.29$	5.73 2.32	6.47 2.43
DKT ADKF-IFT	4.00 3.44 2.41	3.40 3.19 2.12	2.99 2.14	2.98 2.99 2.26	3.85 2.67 1.38	ADKF-IFT	2.12	1.86	1.68	1.74	1.36

(a) Classification (157 tasks).

ADKF-IFT consistently ranks the best in all settings!

Appendix 2: Statistical Comparisons on FS-Mol

Table 5: *p*-values from the two-sided Wilcoxon signed-rank test for statistical comparisons between ADKF-IFT and DKT/DKT+/ADKF. The null hypothesis is that the median of their performance differences on all FS-Mol test tasks is zero. The significance level is set to $\alpha = 0.05$.

na anti-condition may be an Action		Support set size					
Compared models	Task type	16	32	64	128	256	
ADKF-IFT vs DKT	Classification Regression	${f 1.4 imes 10^{-12}\ 8.2 imes 10^{-2}}$	8.1×10^{-14} 9.6×10^{-2}	$egin{array}{llllllllllllllllllllllllllllllllllll$	$\begin{array}{c} \mathbf{1.0\times10^{-8}}\\ \mathbf{7.1\times10^{-5}}\end{array}$	$\begin{array}{c} \mathbf{3.4\times10^{-7}}\\ \mathbf{9.8\times10^{-7}}\end{array}$	
ADKF-IFT vs DKT+	Classification Regression	$\begin{array}{c} \textbf{3.2}\times \textbf{10^{-13}}\\ \textbf{3.2}\times \textbf{10^{-2}} \end{array}$	7.0×10^{-15} 4.2×10^{-1}	$\begin{array}{c} \textbf{2.3}\times \textbf{10^{-13}}\\ \textbf{3.4}\times \textbf{10^{-5}} \end{array}$	$\begin{array}{c} \mathbf{1.2\times10^{-9}}\\ \mathbf{5.2\times10^{-10}}\end{array}$	$\begin{array}{c} 1.6\times10^{-6}\\ 1.2\times10^{-5}\end{array}$	
ADKF-IFT vs ADKF	Classification Regression	$\begin{array}{c} \mathbf{1.7\times10^{-2}}\\ \mathbf{2.8\times10^{-3}}\end{array}$	1.1×10^{-1} 4.2×10^{-4}	4.8×10^{-1} 1.3×10^{-3}	8.3×10^{-1} 4.1×10^{-6}	$\begin{array}{c} 1.6\times10^{-3}\\ 1.3\times10^{-5}\end{array}$	

The improvements of ADKF-IFT over other methods are statistically significant!

Appendix 3: Sub-benchmark Performance on FS-Mol

Table 6: Mean performance with standard errors of top performing methods on FS-Mol test tasks within each sub-benchmark (broken down by EC category) at support set size 64 (the median of all considered support sizes). Note that class 2 is most common in the FS-Mol training set ($\sim 1,500$ training tasks), whereas classes 6 and 7 are least common in the FS-Mol training set (< 50 training tasks each).

FS-Mol sub-benchmark (EC category)			Method				
Class	Description	#tasks	RF	GP-ST	ProtoNet	DKT	ADKF-IFT
1	oxidoreductases	7	0.156 ± 0.044	0.152 ± 0.040	0.137 ± 0.037	0.145 ± 0.040	0.160 ± 0.045
2	kinases	125	0.152 ± 0.009	0.161 ± 0.009	0.285 ± 0.010	0.282 ± 0.010	0.299 ± 0.010
3	hydrolases	20	0.229 ± 0.032	0.230 ± 0.032	0.245 ± 0.034	0.254 ± 0.034	0.262 ± 0.033
4	lysases	2	0.276 ± 0.182	0.284 ± 0.189	0.265 ± 0.211	0.272 ± 0.206	0.279 ± 0.201
5	isomerases	1	0.166 ± 0.040	0.212 ± 0.052	0.172 ± 0.044	0.204 ± 0.058	0.198 ± 0.046
6	ligases	1	0.149 ± 0.035	0.199 ± 0.028	0.170 ± 0.028	0.229 ± 0.013	0.231 ± 0.022
7	translocases	1	0.128 ± 0.039	0.109 ± 0.049	0.099 ± 0.028	0.122 ± 0.022	0.109 ± 0.033
	all enzymes	157	0.163 ± 0.009	0.171 ± 0.009	0.271 ± 0.009	0.271 ± 0.010	0.285 ± 0.010

(a) Classification ($\Delta AUPRC$).

(b) Regression (R_{os}^2) .

FS-Mol sub-benchmark (EC category)			Method					
Class	Description	#tasks	RF	GP-ST	CNP	DKT	ADKF-IFT	
1	oxidoreductases	6	0.108 ± 0.087	0.103 ± 0.076	-0.012 ± 0.011	0.098 ± 0.078	0.116 ± 0.079	
2	kinases	82	0.160 ± 0.019	0.162 ± 0.022	0.127 ± 0.017	0.343 ± 0.022	0.363 ± 0.024	
3	hydrolases	19	0.256 ± 0.058	0.267 ± 0.061	0.014 ± 0.015	0.295 ± 0.063	0.310 ± 0.062	
4	lysases	2	0.418 ± 0.405	0.417 ± 0.416	0.100 ± 0.068	0.440 ± 0.418	0.442 ± 0.403	
5	isomerases	1	0.125 ± 0.077	0.086 ± 0.082	-0.012 ± 0.010	0.209 ± 0.113	0.226 ± 0.063	
6	ligases	1	0.182 ± 0.040	0.202 ± 0.079	0.002 ± 0.004	0.277 ± 0.035	0.279 ± 0.043	
,	all enzymes	111	0.178 ± 0.019	0.181 ± 0.021	0.097 ± 0.014	0.321 ± 0.021	0.340 ± 0.022	

Appendix 4: Meta-testing Cost on FS-Mol



Figure 5: Wall-clock time consumed (with standard errors) when meta-testing on a pre-defined set of FS-Mol classification tasks using each of the compared meta-learning methods.

Appendix 6: Few-shot Molecular Property Prediction on MoleculeNet (Wu et al., 2018)

Table 1: Mean test performance (AUROC%) with standard deviations of all compared methods on MoleculeNet benchmark tasks at support set size 20 (i.e., 2-way 10-shot).

	MoleculeNet benchmark task (#compounds in total)							
Method	Tox21 (8,014)	SIDER (1,427)	MUV (93,127)	ToxCast (8,615)				
Siamese ProtoNet MAML TPN EGNN IterRefLSTM PAR ADKF-IFT	$80.40 \pm 0.35 74.98 \pm 0.32 80.21 \pm 0.24 76.05 \pm 0.24 81.21 \pm 0.16 81.10 \pm 0.17 82.06 \pm 0.12 82.43 \pm 0.60$	$71.10 \pm 4.32 \\64.54 \pm 0.89 \\70.43 \pm 0.76 \\67.84 \pm 0.95 \\72.87 \pm 0.73 \\69.63 \pm 0.31 \\\mathbf{74.68 \pm 0.31} \\67.72 \pm 1.21$	$\begin{array}{c} 59.59 \pm 5.13 \\ 65.88 \pm 4.11 \\ 63.90 \pm 2.28 \\ 65.22 \pm 5.82 \\ 65.20 \pm 2.08 \\ 45.56 \pm 5.12 \\ 66.48 \pm 2.12 \\ \textbf{98.18} \pm \textbf{3.05} \end{array}$	63.70 ± 1.26 66.79 ± 0.85 62.74 ± 1.45 63.65 ± 1.57 - 69.72 ± 1.63 72.07 ± 0.81				
Pre-GNN Meta-MGNN Pre-PAR Pre-ADKF-IFT	$\begin{array}{c} 82.14 \pm 0.08 \\ 82.97 \pm 0.10 \\ 84.93 \pm 0.11 \\ \textbf{86.06} \pm \textbf{0.35} \end{array}$	$73.96 \pm 0.0875.43 \pm 0.2178.08 \pm 0.1670.95 \pm 0.60$	$\begin{array}{c} 67.14 \pm 1.58 \\ 68.99 \pm 1.84 \\ 69.96 \pm 1.37 \\ \textbf{95.74} \pm \textbf{0.37} \end{array}$	73.68 ± 0.74 75.12 ± 0.84 76.22 ± 0.13				