





$$\min_{\theta^1, \theta^2, \dots, \theta^r} \sum_{i=1}^I \mathcal{L}_i(\theta^i, \mathcal{D}_i)$$

How and Where to share parameter?

∴ Conditioning!

- Concatenate
- Additive (conditional bias) > These are actually same thing!  
↳ always same??
- Multi-Head architecture
- Multiplicative (conditional scaling) → more expressive  
~~recall = no~~

이제 다른 task에서 공식을 공유하고 있는가?

attention mechanism! → dot product! (vector 생각)

in this case, dimension summing + element-wise multiplication

Neural Network is ~~not~~ Universal Function Approximator

π... these design decision → 적당히 리식 필요, 문제 의존적

Then, ②. Optimizing the objective

$$\min_{\theta} \sum_{i=1}^I \mathcal{L}_i(\theta, \mathcal{D}_i)$$

$$\hat{\mathcal{L}}(\theta, \mathcal{B}) = \sum_{k \in \mathcal{B}} \mathcal{L}_k(\theta, \mathcal{D}_k^b)$$

basic version

negative transfer ⇒ 공유하는 부분을 줄여라!!

(∵ 아님 문제, 표현 공간 문제)

$$\square + \sum_{t=1}^T \|\theta^t - \theta^{t'}\|$$

overfitting ⇒ 공유하는 부분을 늘려라!!

+ objective form은 마땅히 정해져 존재!

(or weighted)

머리도 휴식이 필요해



Multi-task Learning → Solve multiple tasks  $\mathcal{T}_1, \dots, \mathcal{T}_T$  at once  
 Transfer Learning

$$\min_{\theta} \sum_{t=1}^T \mathcal{L}_t(\theta, \mathcal{D}_t)$$

→ (각 task에 학습된 지식을 전이시키기)

$\mathcal{D}_a$ 를  $\mathcal{D}_b$ 로 transfer

$$\phi \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta, \mathcal{D}^{t'}) \quad \text{Fine-tuning}$$

## Meta-Learning Basic

two ways

Mechanistic View = easy to implement

Probabilistic View = high-level

problem definition!

supervised learning

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}) \quad \text{model parameter}$$

training data

$$= \arg \max_{\phi} \left\{ \log p(\mathcal{D} | \phi) + \log p(\phi) \right\}$$

$$= \arg \max_{\phi} \left\{ \frac{1}{n} \sum \log p(y_i | x_i, \phi) + \log p(\phi) \right\}$$

required labeled data!

Can we incorporate additional data?

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-data}})$$

$$\mathcal{D} = \{ (x_i, y_i) \}_{i=1}^K$$

$$\mathcal{D}_{\text{meta-train}} = \{ \mathcal{D}_1, \dots, \mathcal{D}_n \}$$

$$\mathcal{D}_n = \{ (x_i^n, y_i^n) \}_{i=1}^K$$

⇒ 추가적인 데이터를 주는 것!!

(위에서, 각각 meta (1))

what if we don't to keep  $\mathcal{D}_{\text{meta-train}}$  around forever?

$$\arg \max_{\phi} \log p(\phi | \mathcal{D}, \mathcal{D}_{\text{meta-train}}) \approx \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$$

$$\text{where } \theta^* = \arg \max_{\theta} \log p(\theta | \mathcal{D}_{\text{meta-train}})$$

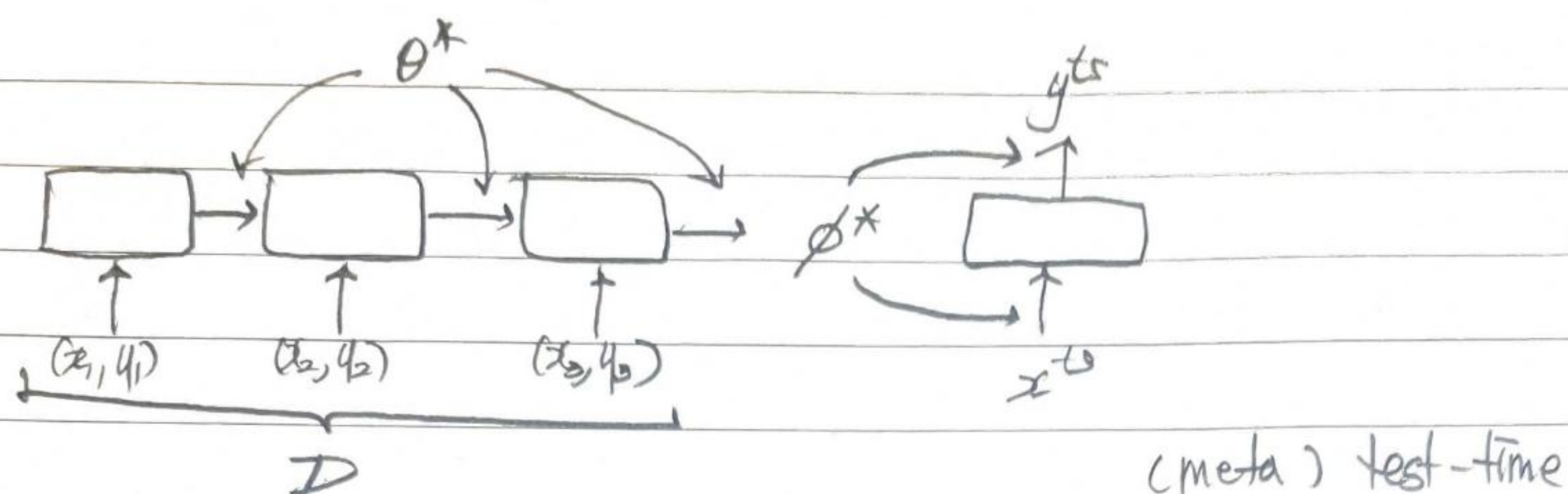
→ meta-learning

$$\phi^* = \arg \max_{\phi} \log p(\phi | \mathcal{D}, \theta^*)$$

→ adaptation

대충원자

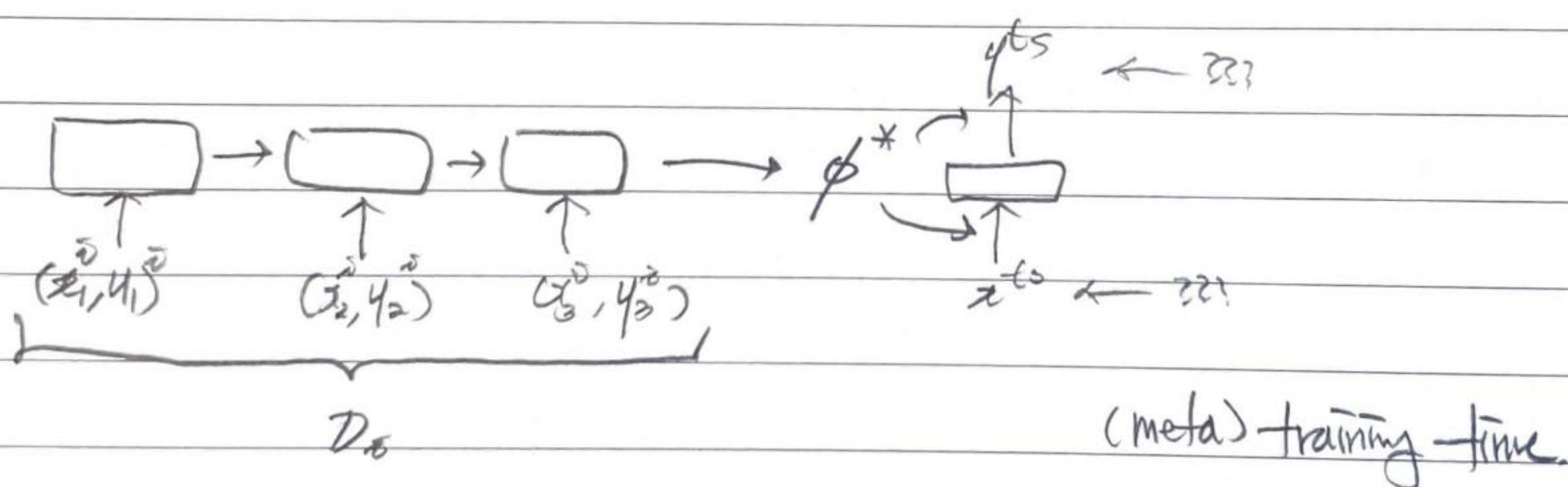




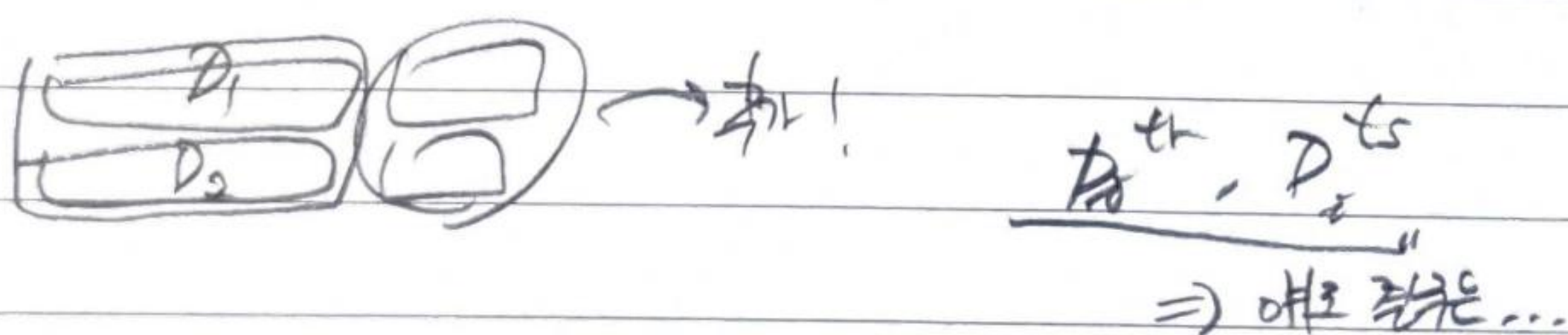
meta learning:  $\theta^* = \arg \max_{\theta} \mathbb{E}_{\mathcal{P}} \mathcal{P}(\phi | \mathcal{D}_{\text{meta-train}})$

adaptation:  $\phi^* = \arg \max_{\phi} \mathcal{P}(\phi | \mathcal{D}_i, \theta^*)$

\* our training procedure is based on a simple machine learning principle:  
: test and train conditions must match.



Reverse a test set for each task!!





\* the complete meta-learning optimization.

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p(\theta | \mathcal{D}_i) \quad \left. \begin{array}{l} \phi_i = \theta \text{ (i.e., } f_{\theta}(\mathcal{D}_i) = \theta) \end{array} \right\}$$

meta learning:  $\theta^* = \underset{\theta}{\operatorname{argmax}} \log p(\theta | \mathcal{D}_{\text{meta-train}})$

multi-task learning

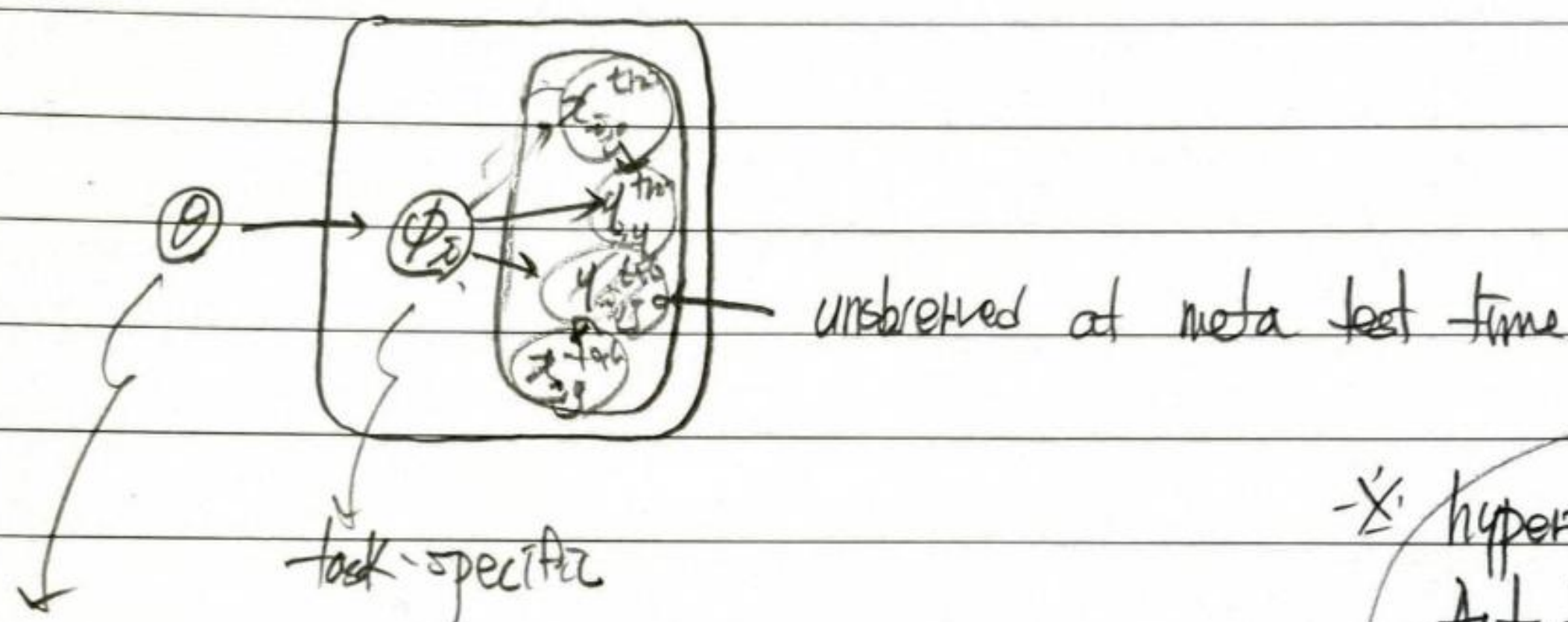
adaptation:  $\phi^* = \underset{\phi}{\operatorname{argmax}} \log p(\phi | \mathcal{D}^{\text{tr}}, \theta^*)$



$$\phi^* = f_{\theta^*}(\mathcal{D}^{\text{tr}})$$

learn  $\theta$  st  $\phi = f_{\theta^*}(\mathcal{D}_i^{\text{tr}})$  is good for  $\mathcal{D}_i^{\text{ts}}$ .

$$\theta^* = \underset{\theta}{\operatorname{max}} \sum_{i=1}^n \log p(\phi_i | \mathcal{D}_i^{\text{ts}}) \quad \text{where } \phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$$



task-specific information!!

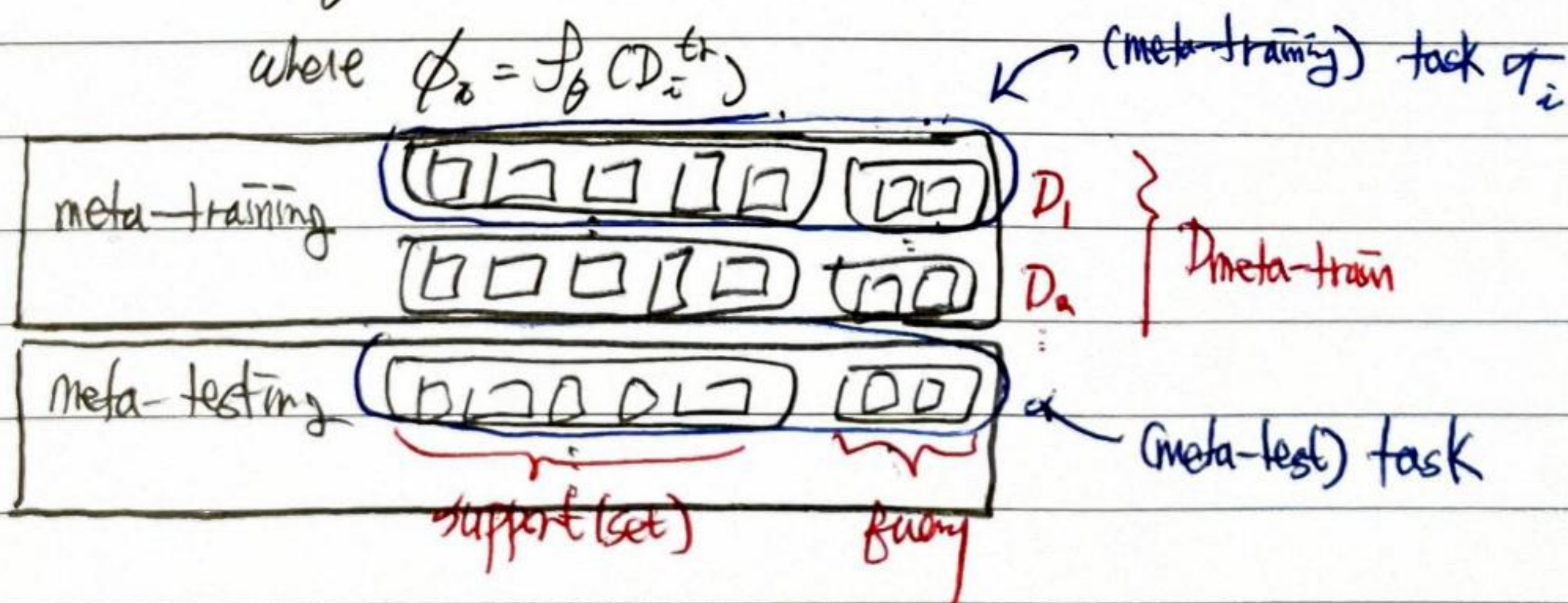
hyperparameter optimization  
Auto-ML  
 $\Rightarrow$  Can be cast as  
meta-learning!!

\* Some meta-learning terminology

learn  $\theta$  st  $\phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$  is good for  $\mathcal{D}_i^{\text{ts}}$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p(\phi_i | \mathcal{D}_i^{\text{ts}})$$

where  $\phi_i = f_{\theta}(\mathcal{D}_i^{\text{tr}})$



$$\mathcal{D}_{\text{meta-train}} = \{ (\mathcal{D}_1^{\text{tr}}, \mathcal{D}_1^{\text{ts}}), \dots, (\mathcal{D}_n^{\text{tr}}, \mathcal{D}_n^{\text{ts}}) \}$$

$$\mathcal{D}_i = \begin{cases} \mathcal{D}_i^{\text{tr}} = \{ (x_1^{\text{tr}}, y_1^{\text{tr}}), \dots, (x_n^{\text{tr}}, y_n^{\text{tr}}) \} \\ \mathcal{D}_i^{\text{ts}} = \{ (x_1^{\text{ts}}, y_1^{\text{ts}}), \dots, (x_n^{\text{ts}}, y_n^{\text{ts}}) \} \end{cases} \rightarrow \text{shot (2e, k-shot, 5-shot)}$$

대충원자