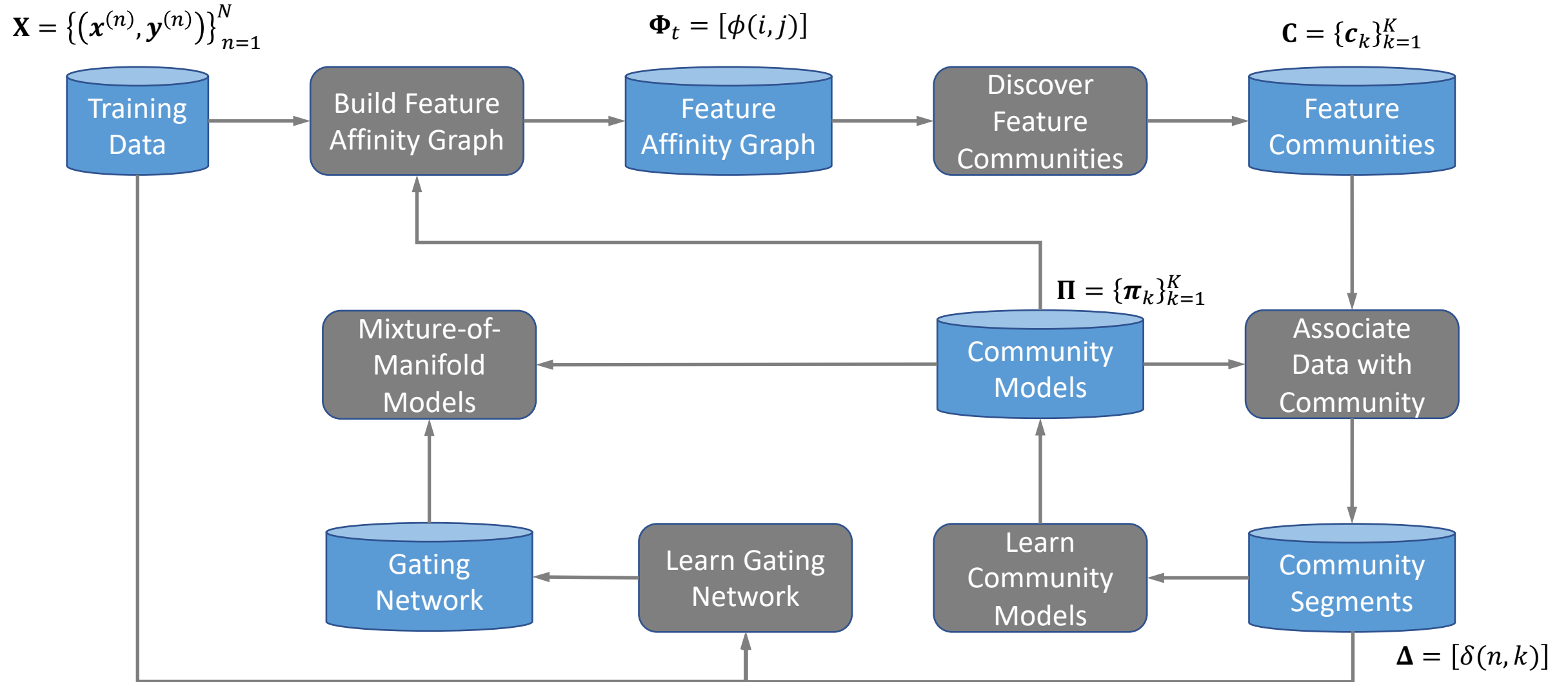


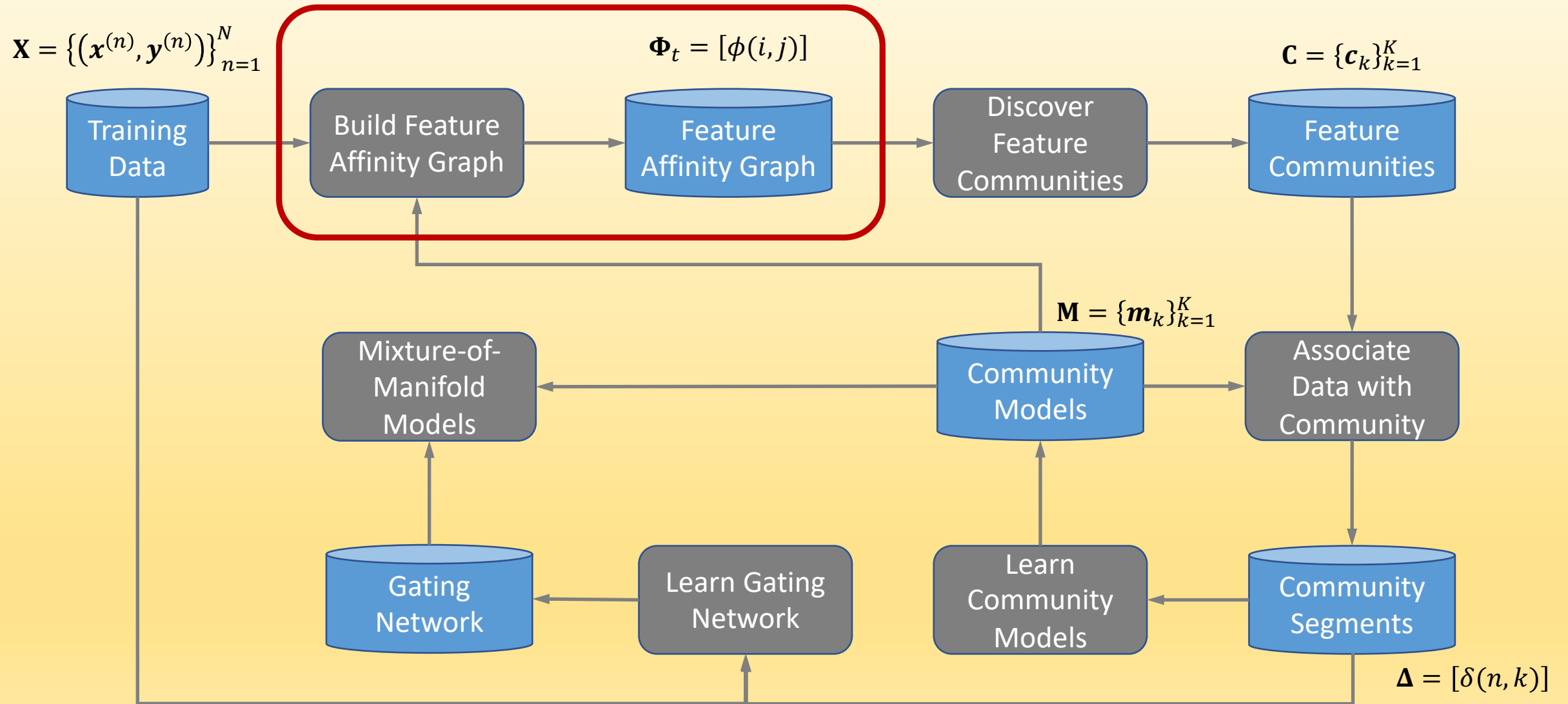
# Mixture-of-Manifold Experts

Shailesh Kumar | Anurag Sahoo

# Overall Architecture



# Build Feature Affinity Graph: $\Phi_0$



# 1. Build Feature Affinity Graph: $\Phi_0$

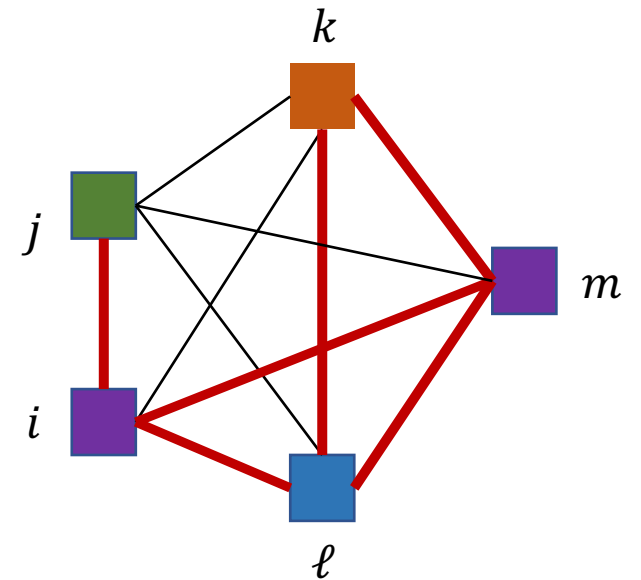
Two features have a high affinity with each other if model accuracy is significantly higher when they are together than independently.

$\theta_0(i)$  = Goodness of model built with feature  $i$  only

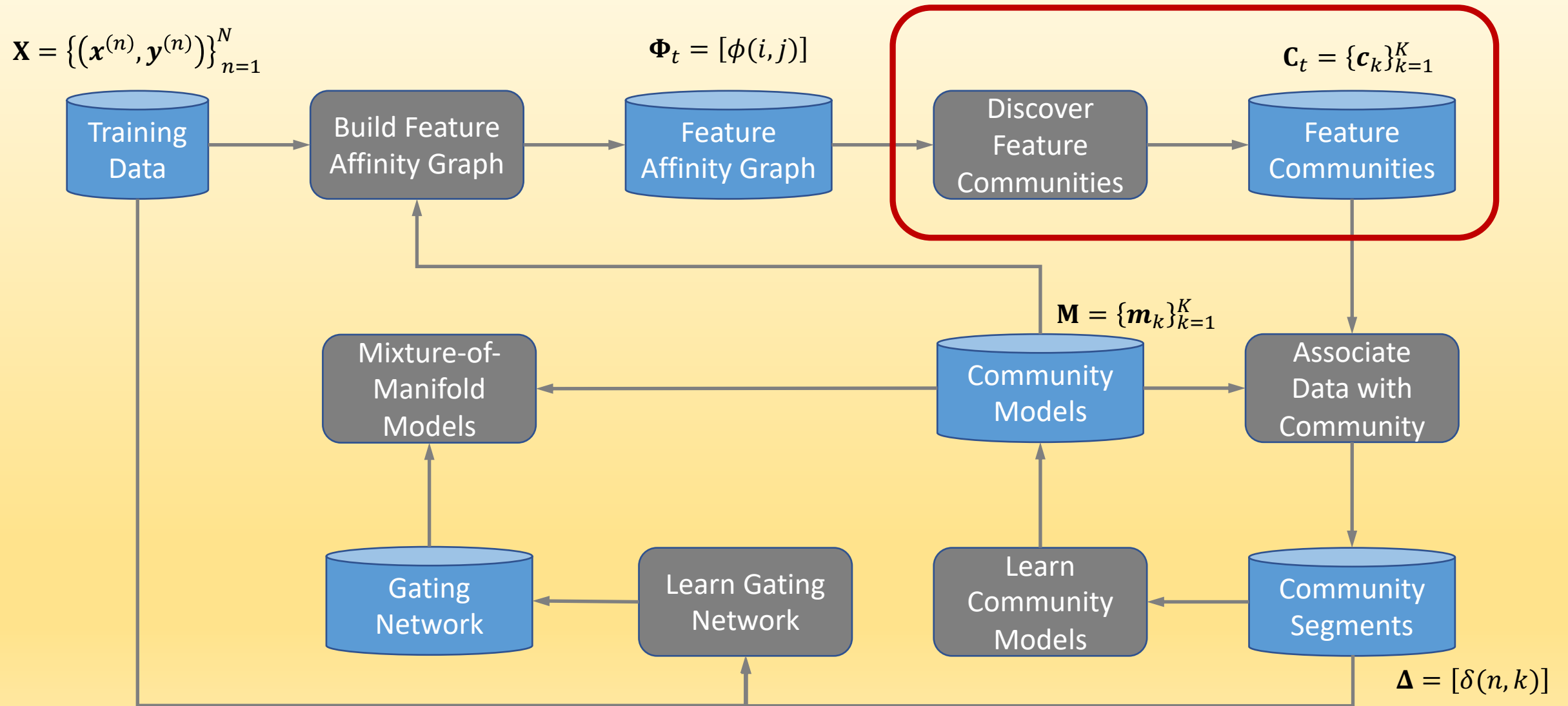
$\theta_0(i, j)$  = Goodness of model built with feature pair  $\{i, j\}$

$\phi_0(i, j)$  = Pair-wise affinity between features  $\{i, j\}$

$$\phi_0(i, j) = \frac{\theta_0(i, j)}{\sqrt{\theta_0(i)\theta_0(j)}}$$



# Discover Feature Communities: $\mathbf{C}_t$



## 2. Discover Feature Communities: $\mathbf{C}_t$

A community (soft maximal clique) of features is such that removing a feature from it or adding a feature to it decreases its coherence.

$\mathbf{c}$  = Any community (set) of features (possibly overlapping with others)

$\pi(\mathbf{c})$  = Coherence of community  $\mathbf{c} = \lambda_1(\mathbf{c}) \times \min\{\mathbf{v}_1(\mathbf{c})\}$

$\lambda_1(\mathbf{c})$  = First Eigen Value of the Affinity Sub-matrix  $\Phi(\mathbf{c})$

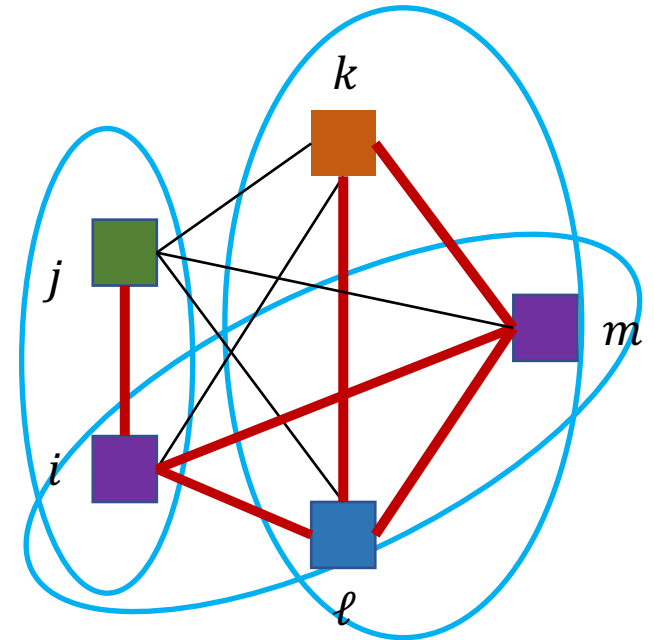
$\mathbf{v}_1(\mathbf{c})$  = First Eigen Vector of the Affinity Sub-matrix  $\Phi(\mathbf{c})$

$\mathbf{c}^*$  = A locally optimal (soft maximal clique) community such that:

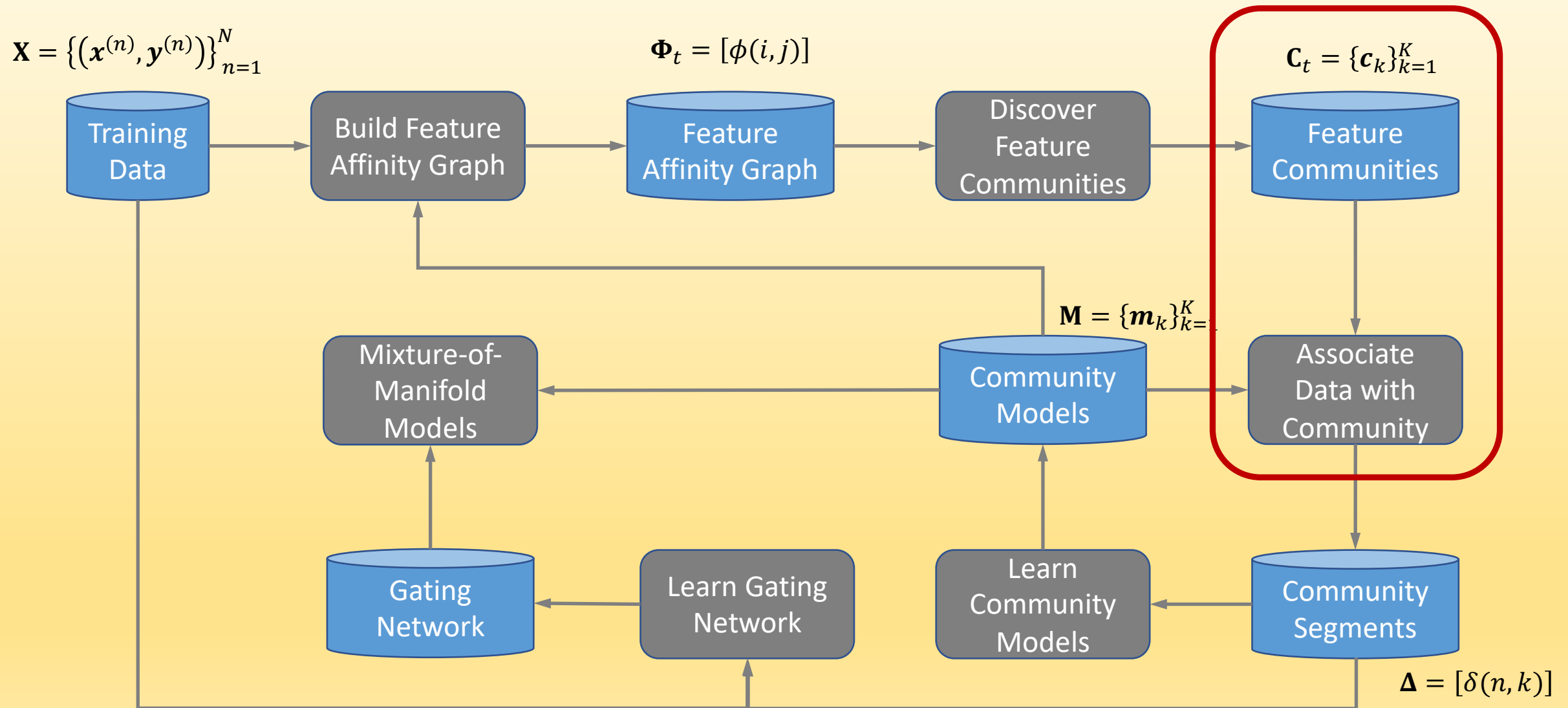
$\pi(\mathbf{c}^*) \geq \pi(\mathbf{c}^* \ominus i), \forall i \in \mathbf{c}^*$  (i.e. removing a feature reduces coherence)

$\pi(\mathbf{c}^*) \geq \pi(\mathbf{c}^* \oplus i), \forall i \notin \mathbf{c}^*$  (i.e. adding a feature also reduces coherence)

$\mathbf{C}_t$  = Set of all soft maximal communities in the feature affinity graph



# Initialize Data Association with Communities



# Initialize Data Association with Communities

Each Feature Community represents a “Sub-Space Manifold”. We now need to associate each data point with one of the sub-space manifold.

$\Delta = [\delta(n, k)]$  = is the association of  $(\mathbf{x}^{(n)}, y^{(n)})$  with community  $\mathbf{c}_k$

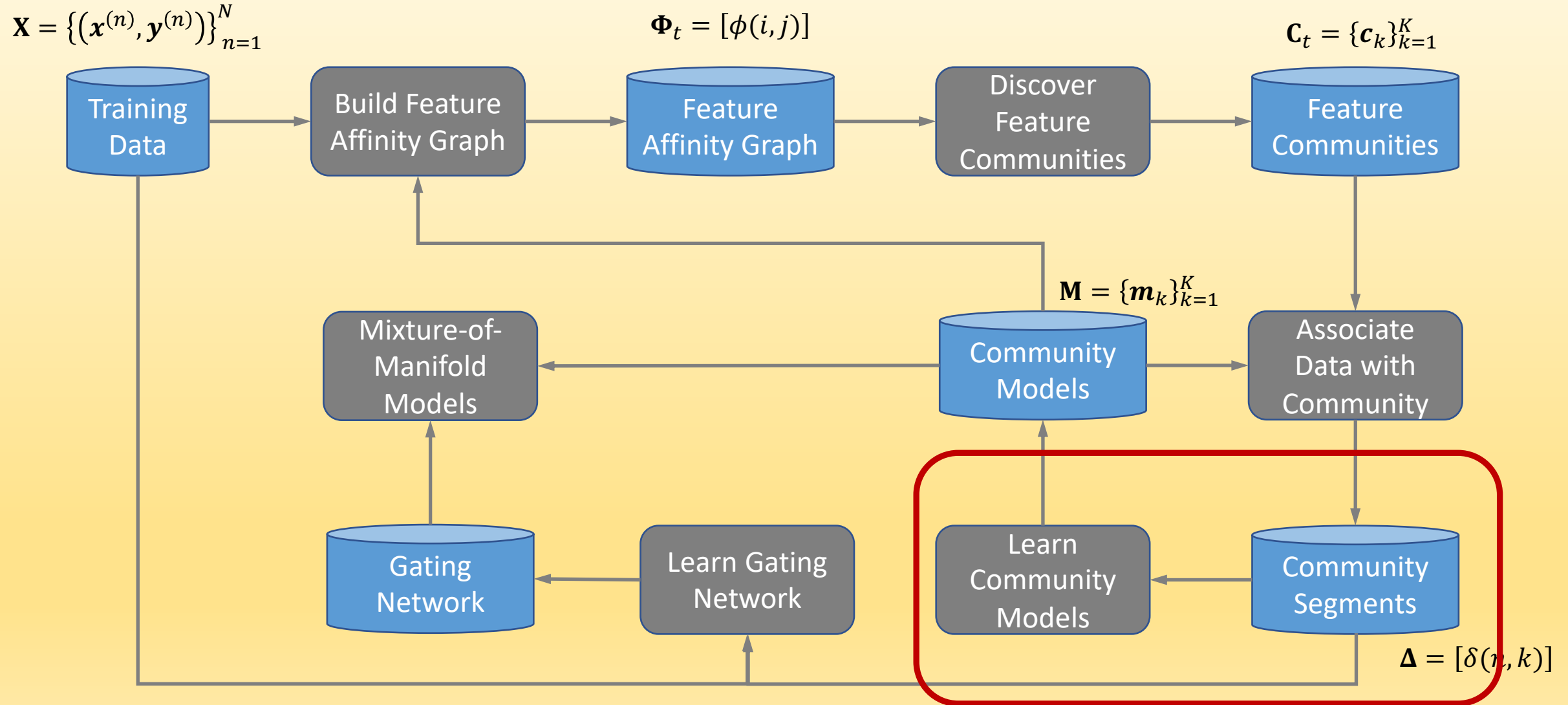
	$\mathbf{c}_1$	...	$\mathbf{c}_k$	...	$\mathbf{c}_K$
$(\mathbf{x}^{(1)}, y^{(1)})$					
$(\mathbf{x}^{(2)}, y^{(2)})$					
$(\mathbf{x}^{(n)}, y^{(n)})$					
$(\mathbf{x}^{(N)}, y^{(N)})$					

$\delta_0(n, k) = \frac{1}{K}$

Initially all data points are equally associated with all communities



# Learn Community Models



# Learn Community Models

Given a modelling technique (expert type), association of a data point with a community, and the community features, learn a model.

$\mathbf{M}_t = \{\mathbf{m}_k^{(t)}\}_{k=1}^K$  = Set of models trained in iteration  $t$  of association loop.

	$c_1$	...	$c_k$	...	$c_K$
$(\mathbf{x}^{(1)}, y^{(1)})$					
$(\mathbf{x}^{(2)}, y^{(2)})$					
$(\mathbf{x}^{(n)}, y^{(n)})$					
$(\mathbf{x}^{(N)}, y^{(N)})$					
	$\delta_t(*, k)$				

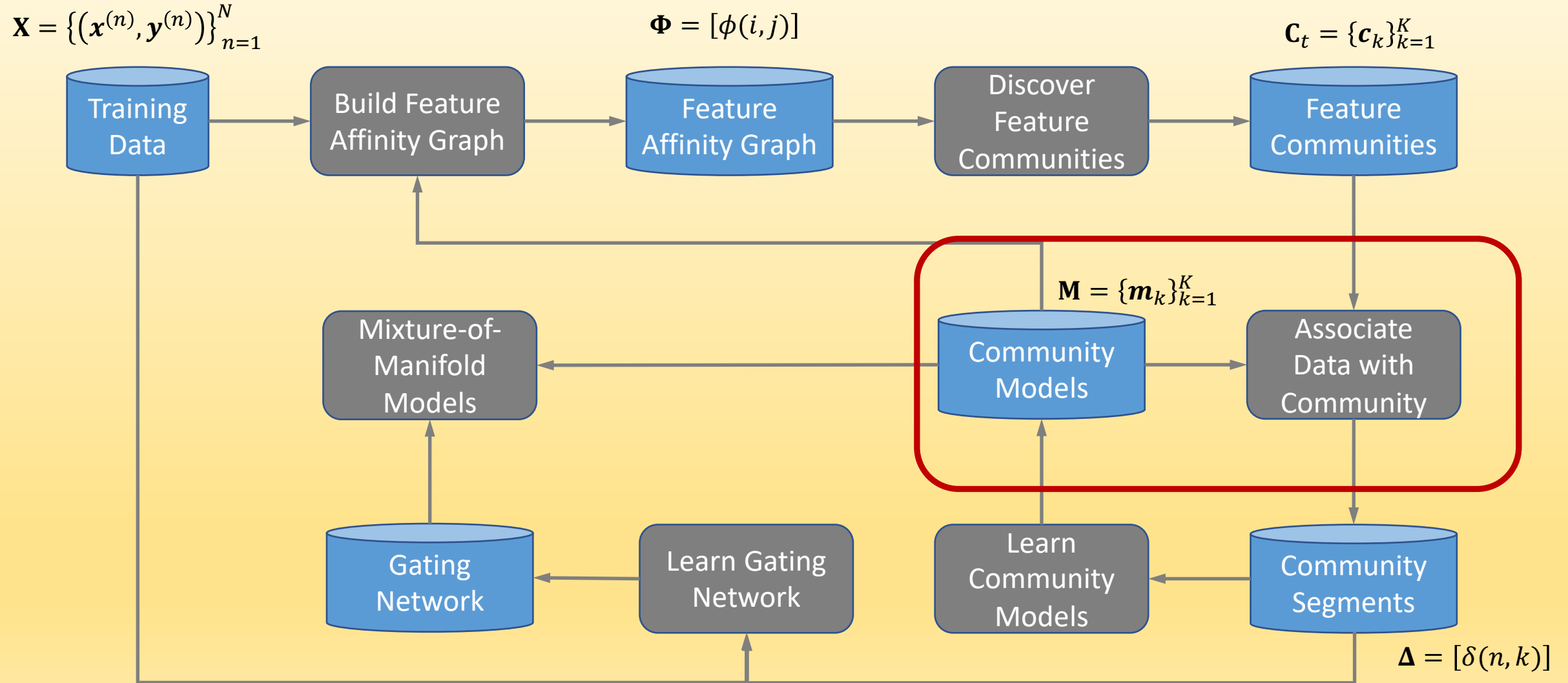
$$\mathbf{m}_k^{(t)} = \text{Train}(\mathbf{X}, \delta_t(*, k)), \forall k = 1 \dots K$$

$\delta_t(n, k)$  = Weight of the  $n$ th data point for  $k$ th model.

In the HARD version, assign each data point to max weighted.

In the SOFT version, assign each data point proportional to weigh.

# Update association of data with community



# Update Association of data with Community

A data point should be associated with that community whose current model gives highest confidence, correct label for this data point

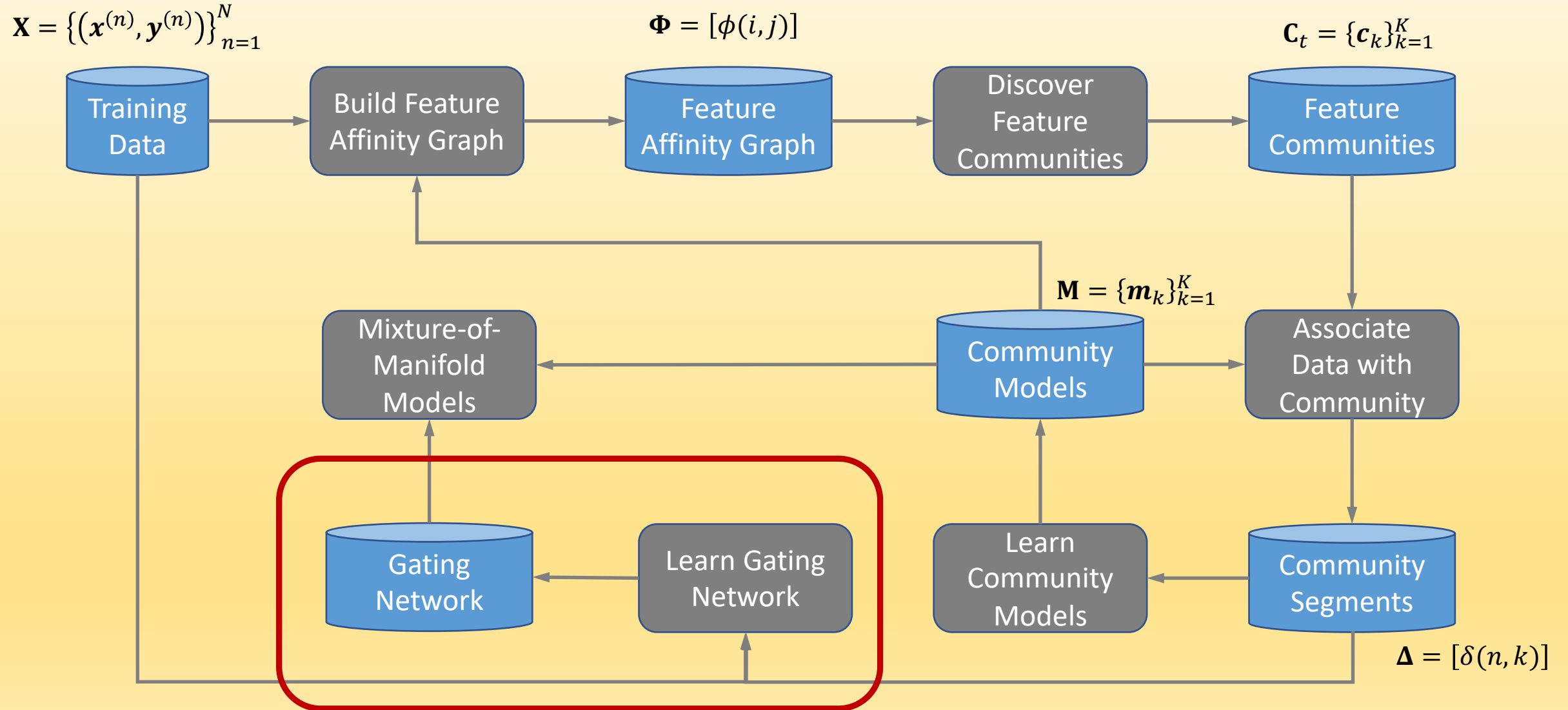
$(\ell_k^{(n)}, f_k^{(n)}) = \text{Predict}(\mathbf{m}_k^{(t)}, \mathbf{x}^{(n)})$  = Predict label and confidence for each data point w.r.t. each model

$\ell_k^{(n)} \in \{-1, 1\}$  class label (two class problem)

$f_k^{(n)} = P(\ell_k^{(n)} | \mathbf{x}_k^{(n)})$ ,  $\forall k = 1 \dots K$  = confidence score in the predicted class label

$\delta_{n,k}^{(t+1)} = \frac{\exp(\beta_t \times y_k^{(n)} \times \ell_k^{(n)} \times f_k^{(n)})}{\sum_j \exp(\beta_t \times y_j^{(n)} \times \ell_j^{(n)} \times f_j^{(n)})}$  = Updated association of  $(\mathbf{x}^{(n)}, y^{(n)})$  with community  $k$

# Learn Gating Network



# Learn Gating Network

When a new data point arrives, the gating network first determines the weightage of each of the experts and combines their outputs.

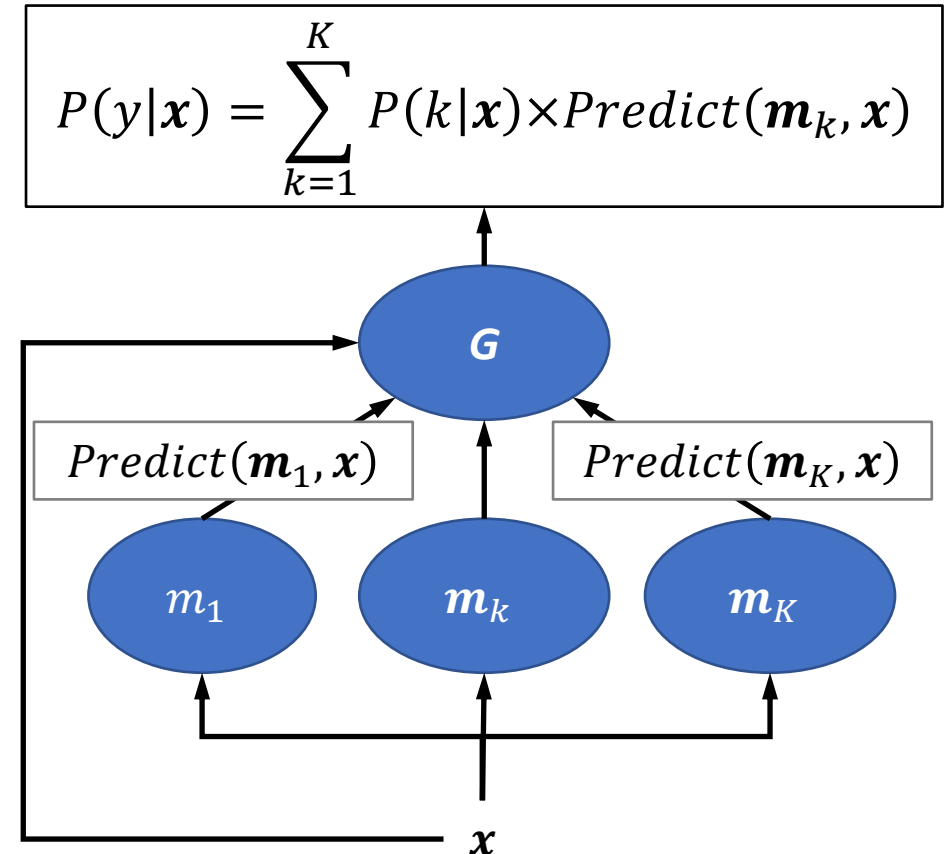
$\mathbf{x}^{(n)}$  = Training Input to the gating network

$\delta(n, k)$  = Training Output of the gating network

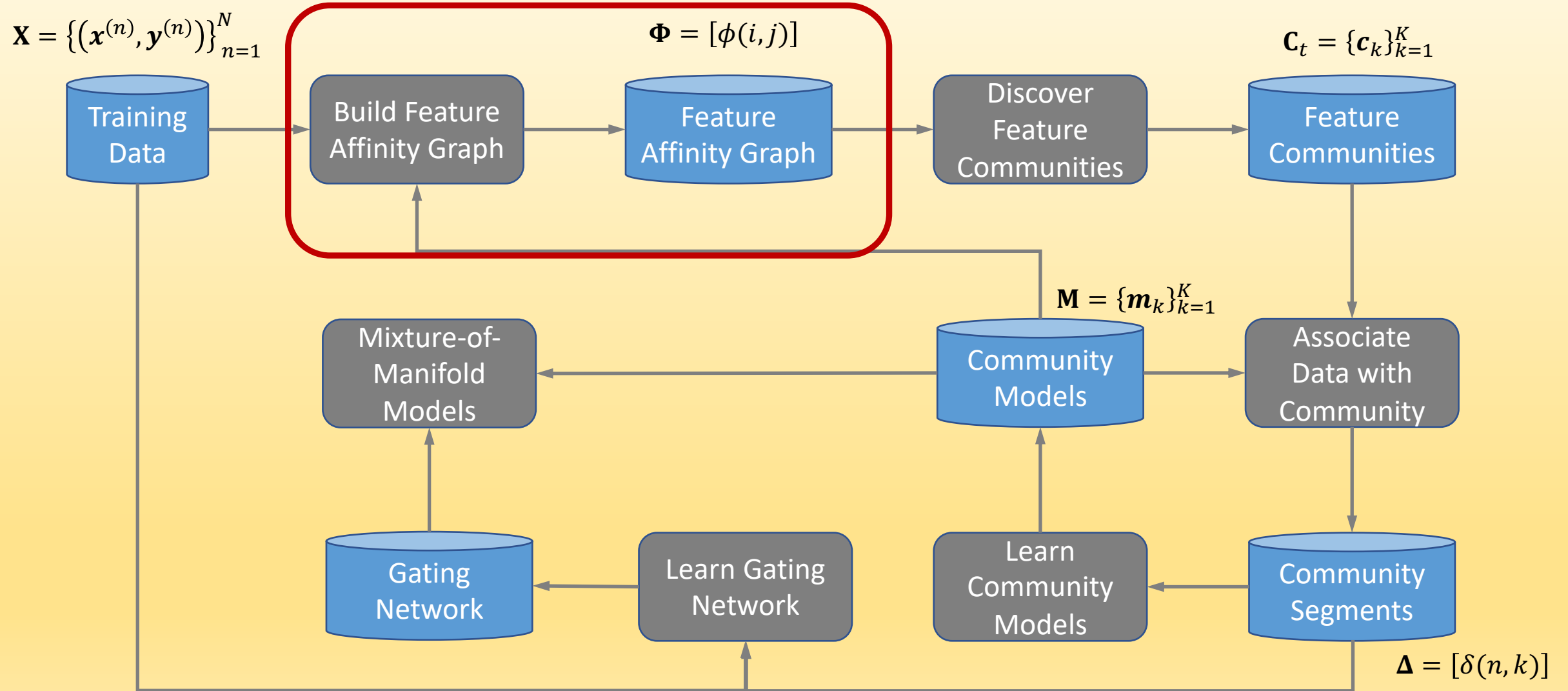
$G(\mathbf{x})$  = Output of the trained gating network

$$G(\mathbf{x}) = (P(1|\mathbf{x}), P(2|\mathbf{x}), \dots, P(K|\mathbf{x}))$$

$P(k|\mathbf{x})$  = probability that gate  $k$  is best for input  $\mathbf{x}$



# Update Feature Affinity Matrix



# Update Feature Affinity Matrix

The initial Affinity Matrix was created based on all data points and only on pair-wise combinations. Now can we do better?

$w(i|k)$  = Importance of feature  $i$  in model  $k$

$w(i|k) = 0$  if feature  $i$  is not present in model  $k$

$A(k)$  = Accuracy/AUC (some metric) of model  $k$

$\theta_t(i)$  = Expected Goodness of feature  $i$

$$\theta_t(i) = \frac{\sum_{k=1}^K w(i|k)A(k)}{\sum_{k=1}^K w(i|k)}$$

$\theta_t(i, j)$  = Expected Goodness of feature pair  $i, j$

$$\theta_t(i, j) = \frac{\sum_{k=1}^K w(i|k)w(j|k)A(k)}{\sum_{k=1}^K w(i|k)w(j|k)}$$

$\phi_0(i, j)$  = Pair-wise affinity of feature pair  $i, j$

$$\phi_t(i, j) = \frac{\theta_t(i, j)}{\sqrt{\theta_t(i)\theta_t(j)}}$$

