$t=0$   (active sampling 전)
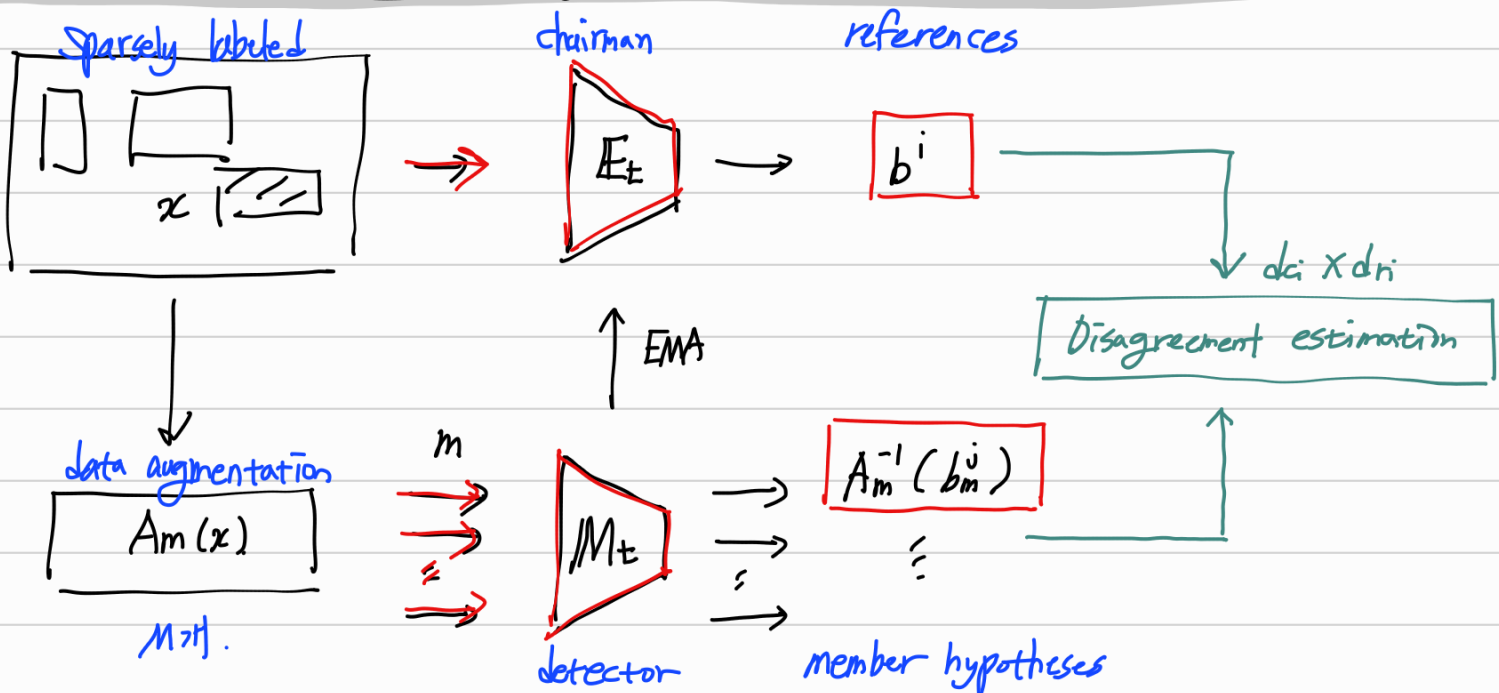
(randomly sampled. fully annotated w bounding boxes)

$$I \longrightarrow L_0$$
$$\phantom{I} \longrightarrow U_0$$

$M(\theta_0)$ : generic object detector

---

$t \geq 1$   ( active sampling 후 )

$$I \longrightarrow L_1 \longrightarrow S_1 \text{ (: sparsely labeled images )}$$
$$\phantom{I} \longrightarrow U_1$$

$M(\theta_t)$   updated   w   labeled images   $L_t \cup S_t$

or   all images   $L_t \cup S_t \cup U_t$

$S_t$ :   sparsely labeled images



Sparsely labeled
$x$

chairman
$E_t$

references
$b^i$

$d_{ci} \times d_{ri}$

Disagreement estimation

EMA

data augmentation
$A_m(x)$

$M > H.$

$m$

$M_t$

detector

$A_m^{-1}(b_m^{ij})$

member hypotheses

chairman : $\mathbb{E}(\theta')$ of $M(\theta)$ : generate box references.
　　　　　　EMA　　　　　　detector

$$\theta_{tr}' = \alpha\, \theta_{tr-1}' + (1-\alpha)\cdot \theta_{tr}$$

($t_{tr}$ : training step within one cycle)

① 
- Disagreement on classification.

1. Given chairman의 예측 box candidates : $\{b^i\}$ >

member boxes $\{b_m^{\ddot{v}}\}$ are assigned to each reference box
　　　　　　　　　　　　　　　　　　　　　　　　　in $\{b^i\}$
using detector-defined assignment strategy > such as max-IoU
　　　　　　　　　　　　　　　　　　　　　　　　　assigner.

2. Given matched pair of boxes $\boxed{\{b^i, b_m^{\ddot{v}}\}}$ >
　　classification disagreement 측정.
　: CE b/w one-hot chairman prediction $g^i$ &
　　　　　posterior predictive member distribution $g^{\ddot{v}}$

$$\boxed{d_c^{ij} = -\mathbb{E}_{g^i}\left[\log g^j\right]}$$

3. Disagreement about box $b^i$ is aggregated among $M$
　　　　　　　　　　　　　　　　　　　　　　Committee members.

$$\boxed{d_c^i = \frac{1}{M}\cdot\sum_m^M \left(\frac{1}{K_{mi}}\sum_j^{K_{mi}} d_c^{ij}\right)}$$

　　　　　　　　　　　　　　　　　　　　　stochastic view.)
( $K_{mi}$ : # of positively matched member predictions in $m$-th

(2)

— ‖ Disagreement on localization ‖

1. $\{b^i, b_m^j\}$ : matched pair of boxes.

2. $\{A_m^{-1}(b_m^j)\}$ : inverse transformations on those boxes.

   fed into localization branch of chairman model $\mathbb{E}^{reg}$.

3. Disagreement over the location of $b^i$ is measured based on chairman re-calibrated boxes.

$$d_r^i = \frac{1}{4} \cdot \sum_k^4 \hat{\sigma}_k \left( \{ \mathbb{E}^{reg}(A_m^{-1}(b_m^j)) \} \right)$$

localization task ⇒ 4-way regression task based on coordinates.

— ‖ Our scoring function for box-level detection task ‖

$$d^i = d_c^i \times d_r^i$$

— Supervised loss for labeled images

· fully labeled images. $\{x_\ell^i\}$ from $L_\pm$ ; $N_\ell > 1$

$$\underset{\text{classification loss}}{\mathcal{L}_\ell} = \frac{1}{N_\ell} \sum_i \mathcal{L}_{cls}(x_\ell^i, y_\ell^i) + \mathcal{L}_{loc}(x_\ell^i, t_\ell^i)$$

classification loss → gt class label ← localization loss function

fully labeled image 개수. fully labeled image. corresponding box location.

— Pseudo - active synergy for sparse images

sparse gt label

$$G(y_s, \hat{y}_{sc})$$    pseudo label

$$= y_s \cup \{ \hat{y}_{sc}^i \mid IOU( \hat{b}_{sc}^i, b_s^j ) \leq \lambda_g, \forall b_s^j \in b_s \}$$

corresponding box

↳ jaccard overlap

$$G( t_s, \hat{t}_{sr} ) \quad \Sigma \quad 일치하지.$$

∴ supervision quality for sparse images.

$$\mathcal{L}_s = \frac{1}{N_c} \cdot \sum_i^{N_s} \mathcal{L}_{cls}( A(x_s^i), G(y_s^i, \hat{y}_{sc}^i) ) + $$
$$\mathcal{L}_{loc}( A(x_s^i), G( t_s^i, \hat{t}_{sr}^i ) )$$

sparsely labeled image 개수.