

# Minimax Classification under Concept Drift with Multidimensional Adaptation and Performance Guarantees

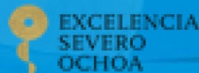
Basque Center for Applied Mathematics-BCAM

Verónica Álvarez, [valvarez@bcamath.org](mailto:valvarez@bcamath.org)

Santiago Mazuelas, [smazuelas@bcamath.org](mailto:smazuelas@bcamath.org)

Jose A. Lozano, [jlozano@bcamath.org](mailto:jlozano@bcamath.org)

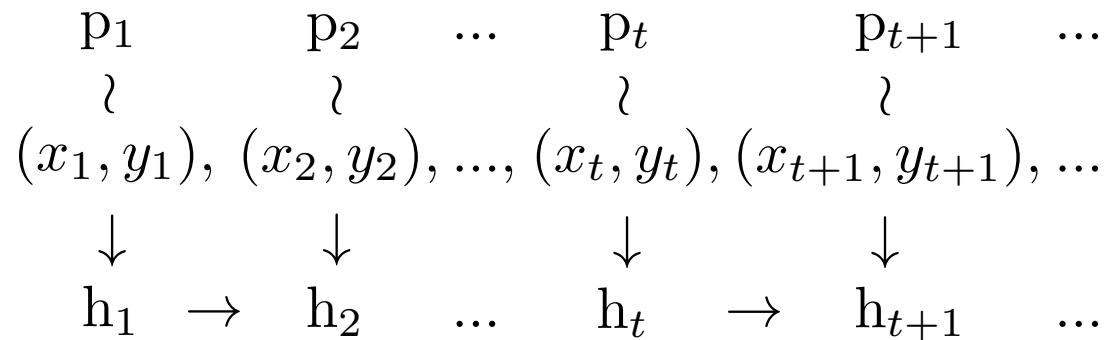
(bcam)



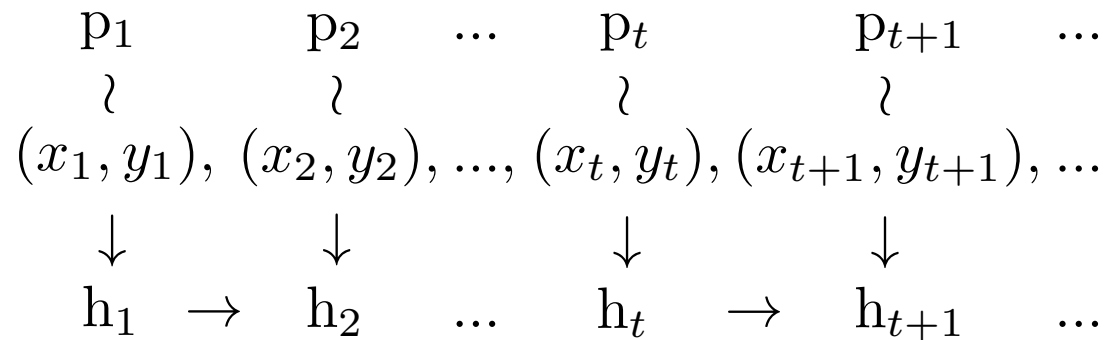
[www.bcamath.org](http://www.bcamath.org)  
basque center for applied mathematics



# Supervised classification under concept drift



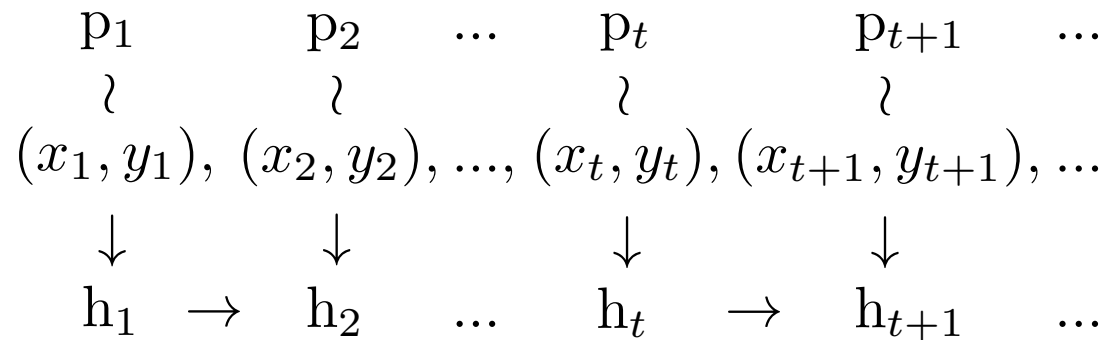
# Supervised classification under concept drift



For instance,

$$h_{t+1} = h_t$$

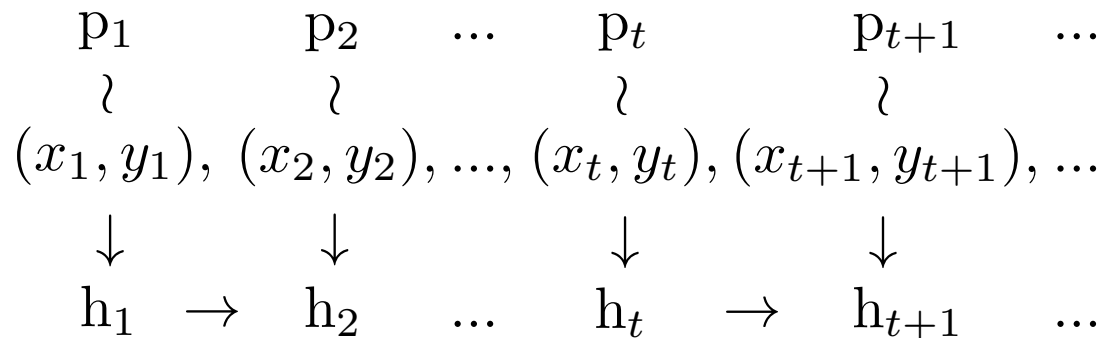
# Supervised classification under concept drift



For instance,

$$h_{t+1} = h_t - k \nabla \ell(h_t, (x_t, y_t))$$

# Supervised classification under concept drift

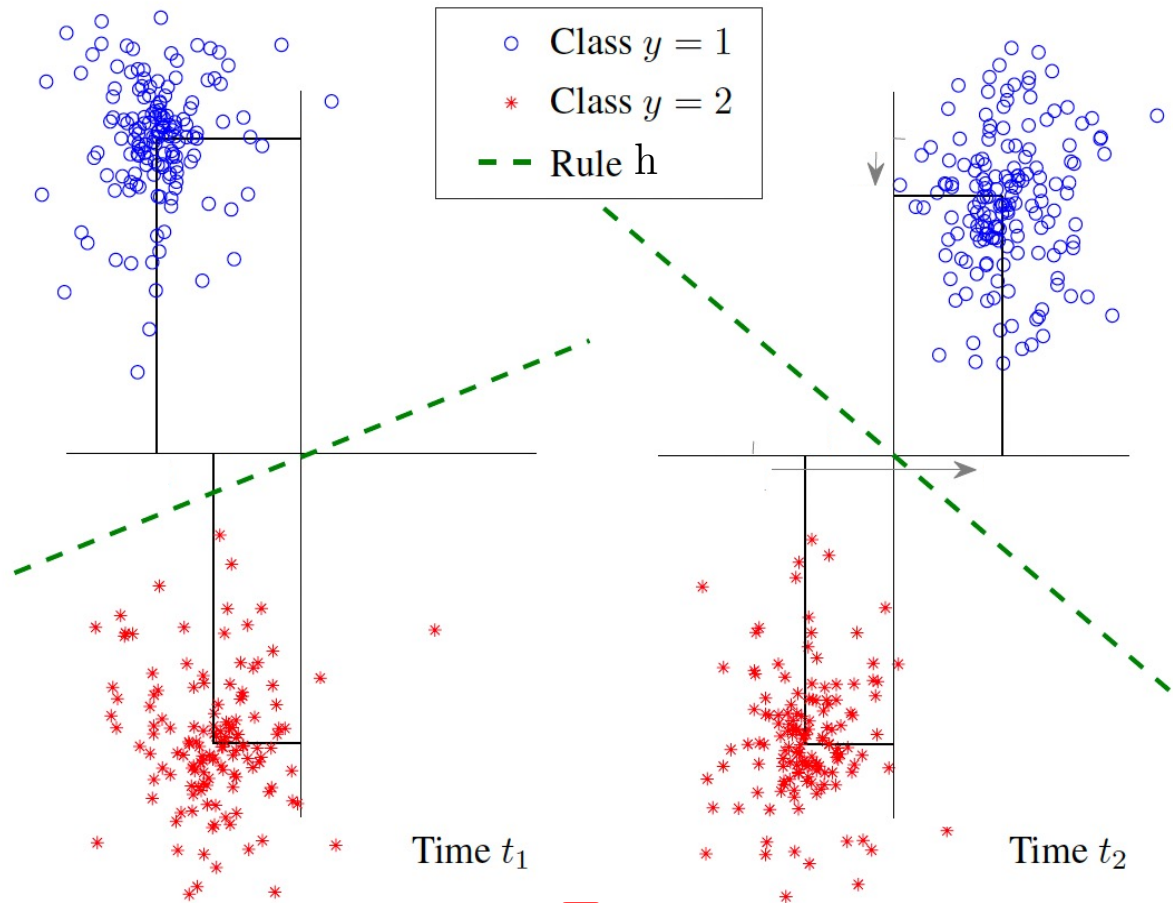


For instance,

$$h_{t+1} = h_t - k \nabla \ell(h_t, (x_t, y_t))$$

$k \in \mathbb{R}$  accounts for a global rate of change

# Supervised classification under concept drift



$$h_{t+1} = h_t - k \nabla \ell(h_t, (x_t, y_t))$$

$k \in \mathbb{R}$  accounts for a global rate of change

# Knowledge gaps

Account for a scalar rate of change: learning rate or forgetting factor

Provide bounds in terms of non-computable quantities: discrepancies between consecutive distributions

# Key contributions

Account for multidimensional time changes: model the evolution of each statistical characteristic of instance-label pairs

Provide computable tight bounds for: error probabilities and accumulated mistakes

# Methodology of Adaptive Minimax Risk Classifiers

Multidimensional adaptation


$$\tau_t = \mathbb{E}_{p_t} \{ \Phi(x, y) \}$$

Feature mapping

$$\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$$

Uncertainty set

$$\mathcal{U}_t = \{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p \{ \Phi(x, y) \} - \hat{\tau}_t| \preceq \lambda_t \}$$

$$\mathcal{U}_t$$

$$\begin{array}{c} \hat{\tau}_t \\ \lambda_t \end{array}$$



# Methodology of Adaptive Minimax Risk Classifiers

Multidimensional adaptation

$$\boldsymbol{\tau}_t = \mathbb{E}_{p_t} \{ \Phi(x, y) \}$$

Feature mapping

$$\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$$

Uncertainty set

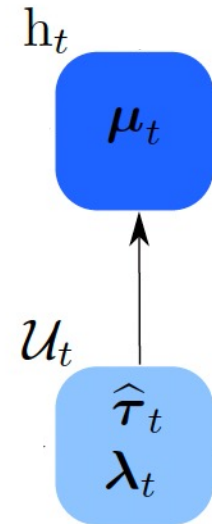
$$\mathcal{U}_t = \{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p \{ \Phi(x, y) \} - \hat{\boldsymbol{\tau}}_t| \leq \boldsymbol{\lambda}_t \}$$

Learning

$$\begin{aligned} & \min_{h \in \mathcal{T}(\mathcal{X}, \mathcal{Y})} \max_{p \in \mathcal{U}_t} \ell(h, p) \\ & = \min_{\boldsymbol{\mu}} 1 - \boldsymbol{\tau}_t^T \boldsymbol{\mu} + \varphi(\boldsymbol{\mu}) + \boldsymbol{\lambda}_t^T |\boldsymbol{\mu}| \end{aligned}$$

Prediction

$$\hat{y} \in \arg \max_{y \in \mathcal{Y}} \Phi(x, y)^T \boldsymbol{\mu}^*$$



# Methodology of Adaptive Minimax Risk Classifiers

Multidimensional adaptation

$$\boldsymbol{\tau}_t = \mathbb{E}_{p_t} \{ \Phi(x, y) \}$$

Feature mapping

$$\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$$

Uncertainty set

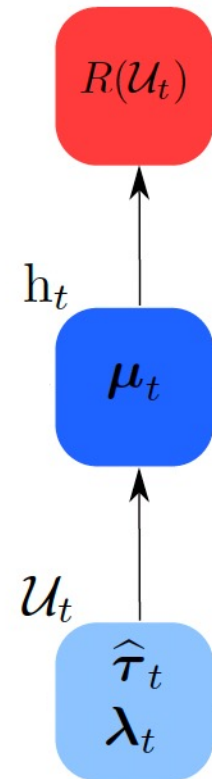
$$\mathcal{U}_t = \{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p \{ \Phi(x, y) \} - \hat{\boldsymbol{\tau}}_t| \leq \boldsymbol{\lambda}_t \}$$

Learning

$$\begin{aligned} & \min_{h \in \mathcal{T}(\mathcal{X}, \mathcal{Y})} \max_{p \in \mathcal{U}_t} \ell(h, p) \\ & = \min_{\boldsymbol{\mu}} 1 - \boldsymbol{\tau}_t^T \boldsymbol{\mu} + \varphi(\boldsymbol{\mu}) + \boldsymbol{\lambda}_t^T |\boldsymbol{\mu}| = R(\mathcal{U}_t) \end{aligned}$$

Prediction

$$\hat{y} \in \arg \max_{y \in \mathcal{Y}} \Phi(x, y)^T \boldsymbol{\mu}^*$$



# Methodology of Adaptive Minimax Risk Classifiers

Multidimensional adaptation

$$\boldsymbol{\tau}_t = \mathbb{E}_{p_t} \{ \Phi(x, y) \}$$

Feature mapping

$$\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^m$$

Uncertainty set

$$\mathcal{U}_t = \{ p \in \Delta(\mathcal{X} \times \mathcal{Y}) : |\mathbb{E}_p \{ \Phi(x, y) \} - \hat{\boldsymbol{\tau}}_t| \preceq \boldsymbol{\lambda}_t \}$$

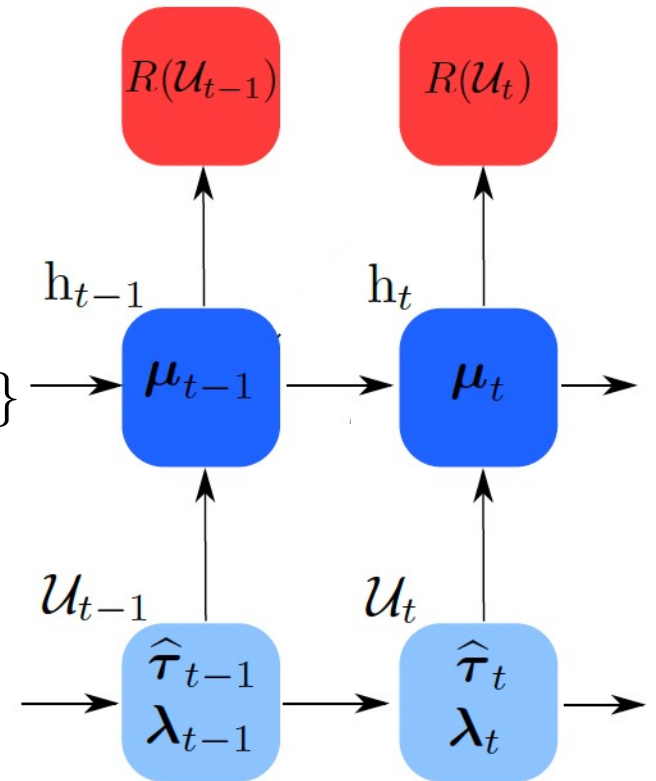
Learning

$$\min_{h \in \mathcal{T}(\mathcal{X}, \mathcal{Y})} \max_{p \in \mathcal{U}_t} \ell(h, p)$$

$$= \min_{\boldsymbol{\mu}} 1 - \boldsymbol{\tau}_t^T \boldsymbol{\mu} + \varphi(\boldsymbol{\mu}) + \boldsymbol{\lambda}_t^T |\boldsymbol{\mu}| = R(\mathcal{U}_t)$$

Prediction

$$\hat{y} \in \arg \max_{y \in \mathcal{Y}} \Phi(x, y)^T \boldsymbol{\mu}^*$$



# Multidimensional Adaptation

Tracking underlying distribution: obtain  $\hat{\tau}_t, \lambda_t$

Dynamical system: model the evolution of each component of  $\tau_t$

$$\eta_{t,i} = \mathbf{H}_t \eta_{t-1,i} + \mathbf{w}_{t,i} \quad \text{with } \eta_{t,i} = [\tau_{t,i}, \tau'_{t,i}, \tau''_{t,i}, \dots, \tau_{t,i}^{(k)}]^T$$
$$\Phi_i(x_t, y_t) = \tau_{t,i} + v_{t,i}$$

# Multidimensional Adaptation

Tracking underlying distribution: obtain  $\hat{\tau}_t, \lambda_t$

Dynamical system: model the evolution of each component of  $\tau_t$

$$\begin{aligned} \eta_{t,i} &= \mathbf{H}_t \eta_{t-1,i} + \mathbf{w}_{t,i} \\ \Phi_i(x_t, y_t) &= \tau_{t,i} + v_{t,i} \end{aligned} \quad \text{with } \eta_{t,i} = \left[ \tau_{t,i}, \tau'_{t,i}, \tau''_{t,i}, \dots, \tau_{t,i}^{(k)} \right]^T$$

Unbiased linear estimator of  $\eta_{t,i}$  with minimum MSE

$$\hat{\eta}_{t,i} = \mathbf{H}_t \hat{\eta}_{t-1,i} - \mathbf{k}_{t,i} (\hat{\tau}_{t-1,i} - \Phi_i(x_{t-1}, y_{t-1}))$$

# Multidimensional Adaptation

Tracking underlying distribution: obtain  $\hat{\tau}_t, \lambda_t$

Dynamical system: model the evolution of each component of  $\tau_t$

$$\eta_{t,i} = \mathbf{H}_t \eta_{t-1,i} + \mathbf{w}_{t,i} \quad \text{with } \eta_{t,i} = [\tau_{t,i}, \tau'_{t,i}, \tau''_{t,i}, \dots, \tau_{t,i}^{(k)}]^T$$
$$\Phi_i(x_t, y_t) = \tau_{t,i} + v_{t,i}$$

Unbiased linear estimator of  $\eta_{t,i}$  with minimum MSE

$$\hat{\eta}_{t,i} = \mathbf{H}_t \hat{\eta}_{t-1,i} - \mathbf{k}_{t,i} (\hat{\tau}_{t-1,i} - \Phi_i(x_{t-1}, y_{t-1}))$$

$k_{t,1}, k_{t,2}, \dots, k_{t,m}$  accounts for multidimensional time changes

# Multidimensional Adaptation

Tracking underlying distribution: obtain  $\hat{\tau}_t, \lambda_t$

Dynamical system: model the evolution of each component of  $\tau_t$

$$\eta_{t,i} = \mathbf{H}_t \eta_{t-1,i} + \mathbf{w}_{t,i} \quad \text{with } \eta_{t,i} = [\tau_{t,i}, \tau'_{t,i}, \tau''_{t,i}, \dots, \tau_{t,i}^{(k)}]^T$$
$$\Phi_i(x_t, y_t) = \tau_{t,i} + v_{t,i}$$

Unbiased linear estimator of  $\eta_{t,i}$  with minimum MSE

$$\hat{\eta}_{t,i} = \mathbf{H}_t \hat{\eta}_{t-1,i} - \mathbf{k}_{t,i} (\hat{\tau}_{t-1,i} - \Phi_i(x_{t-1}, y_{t-1}))$$

$k_{t,1}, k_{t,2}, \dots, k_{t,m}$  accounts for multidimensional time changes

## Performance guarantees

Error probability

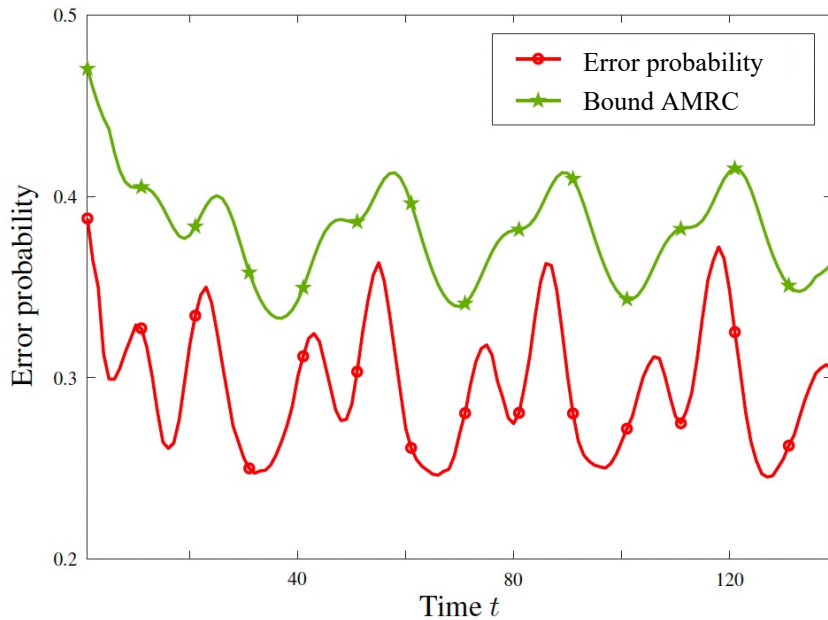
$$R(\mathbf{h}_t) \leq R(\mathcal{U}_t)$$

Accumulated mistakes

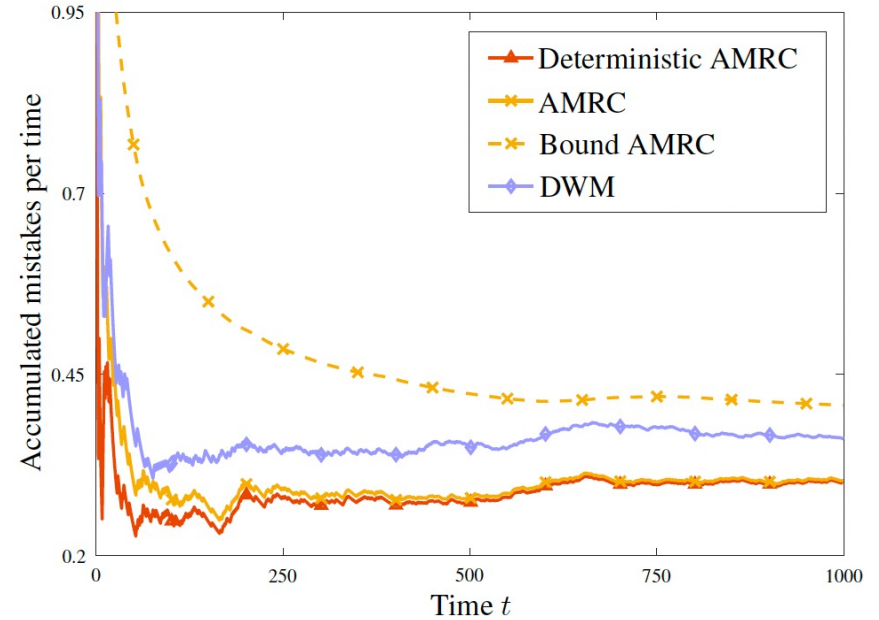
$$\sum_{t=1}^T \mathbb{I}\{\hat{y}_t \neq y_t\} \leq \sum_{t=1}^T R(\mathcal{U}_t) + \sqrt{2T \log \frac{1}{\delta}}$$

# Experimental results

Instantaneous bounds for error probabilities



Bounds for accumulated mistakes



Algorithm	Weather	Elec2	German	Chess	Usenet1	Email	Poker
DWM	<b>30.0</b>	36.3	41.2	35.2	36.3	39.5	22.0
Projectron	30.7	36.8	40.0	35.1	49.1	48.2	22.6
Unid. AMRC	31.3	40.1	30.3	38.3	46.3	48.4	39.4
AMRC	32.3	35.8	30.3	<b>27.7</b>	35.7	43.7	<b>21.9</b>
Det. AMRC	<b>30.0</b>	<b>33.9</b>	<b>30.0</b>	33.4	<b>32.0</b>	<b>33.9</b>	<b>21.9</b>



**Paper**



**Code**

