

# Put your **heatmaps** on! Let's find the facial biomarkers!

Yusuf B. Tanrıverdi, Federico Sukno, Gemma Piella

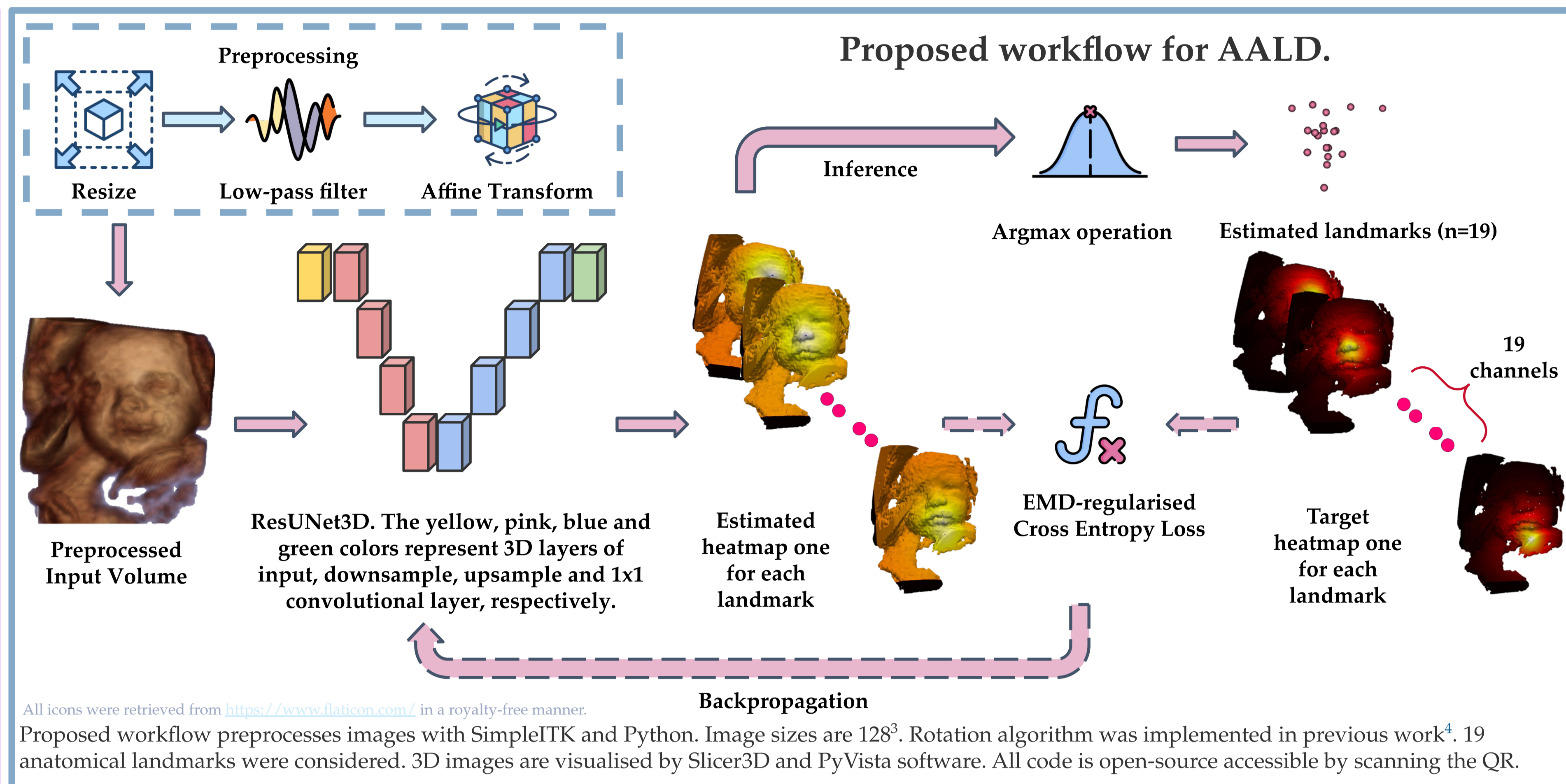


Towards an AI-assisted 3D fetal face reconstruction for craniofacial assessment: **EMD -regularised heatmap regression** for automated anatomical landmark detection (AALD)

## Landmarking leads to a standardised analysis!<sup>1</sup>

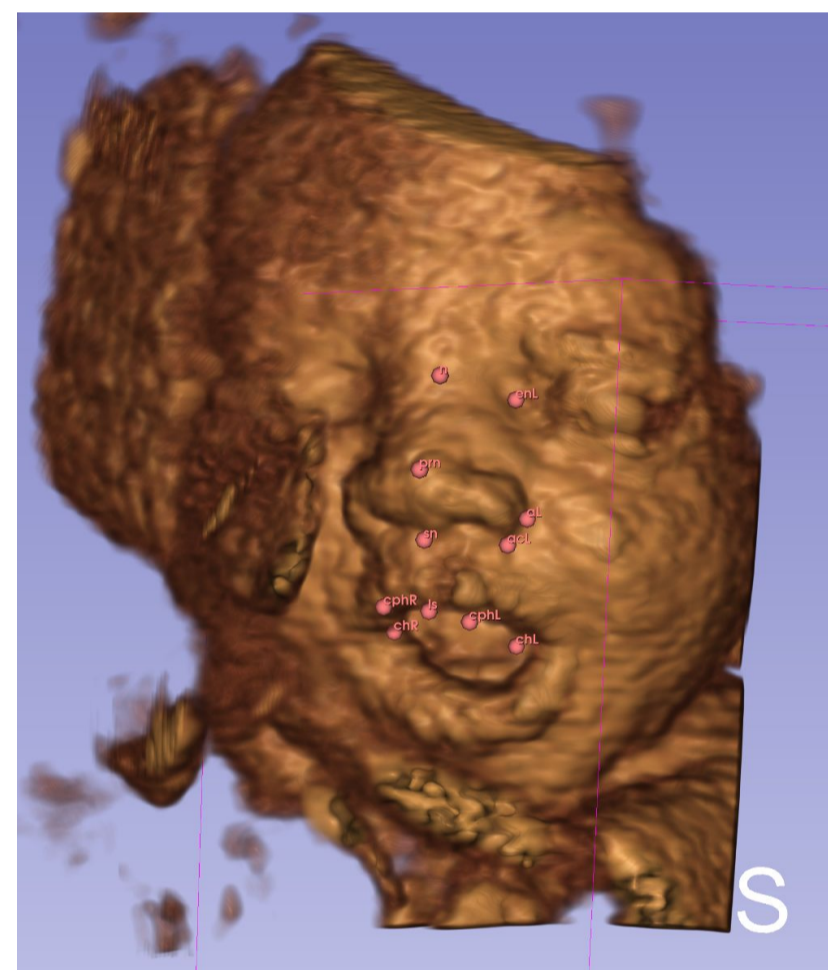
Evaluating fetal facial morphology remains observer-dependent, even with the recent 3D ultrasound trend<sup>2</sup>.

Automated anatomical landmark detection is essential to overcome characteristic limitations of 3D fetal ultrasound, as it serves as the key step toward objective, standardized analysis.



Effective landmark detection requires high-precision probability peaks for every anatomical landmark.<sup>5</sup>

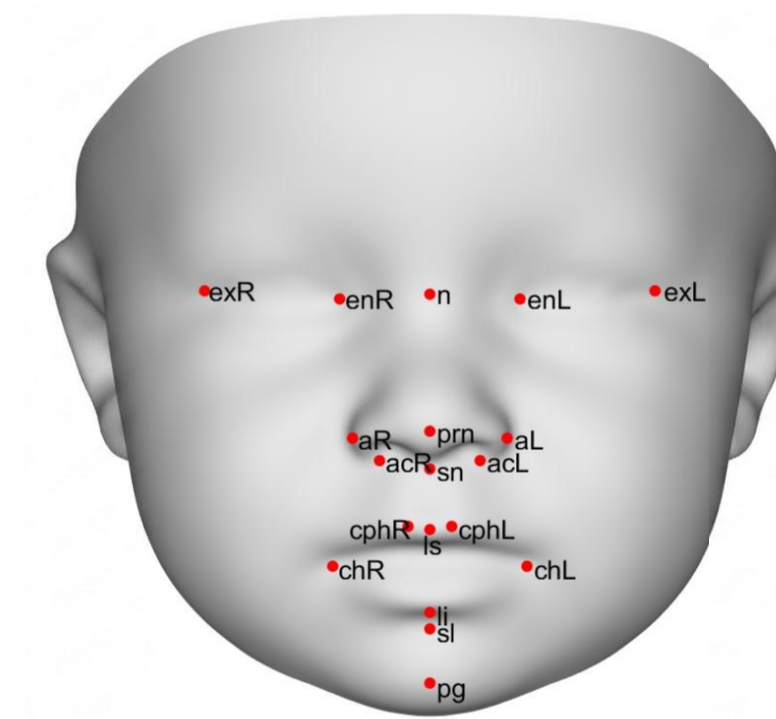
## 3D fetal ultrasound has **limited** acoustic views and facial features might be occluded<sup>3</sup>.



- Thus, a generalisable landmark detector should consider bias and ensure impartiality in the data.
- We collected data from different hospitals in Barcelona: Hospital del Mar, and Hospital Universitari Dexeus. This amounted to 726 images.
- Moreover, our data are spread across gestational ages (GAs), covering a broad range of embryonic development.

## How do we generate ground truth?

Target landmarks are annotated by three different researchers. Invisible facial markers were completed by symmetry and visual comparison.



19 landmarks shown on perinatal face morphology. For each landmark  $(d_0, h_0, w_0)$ , a 3D Gaussian heatmap is generated:

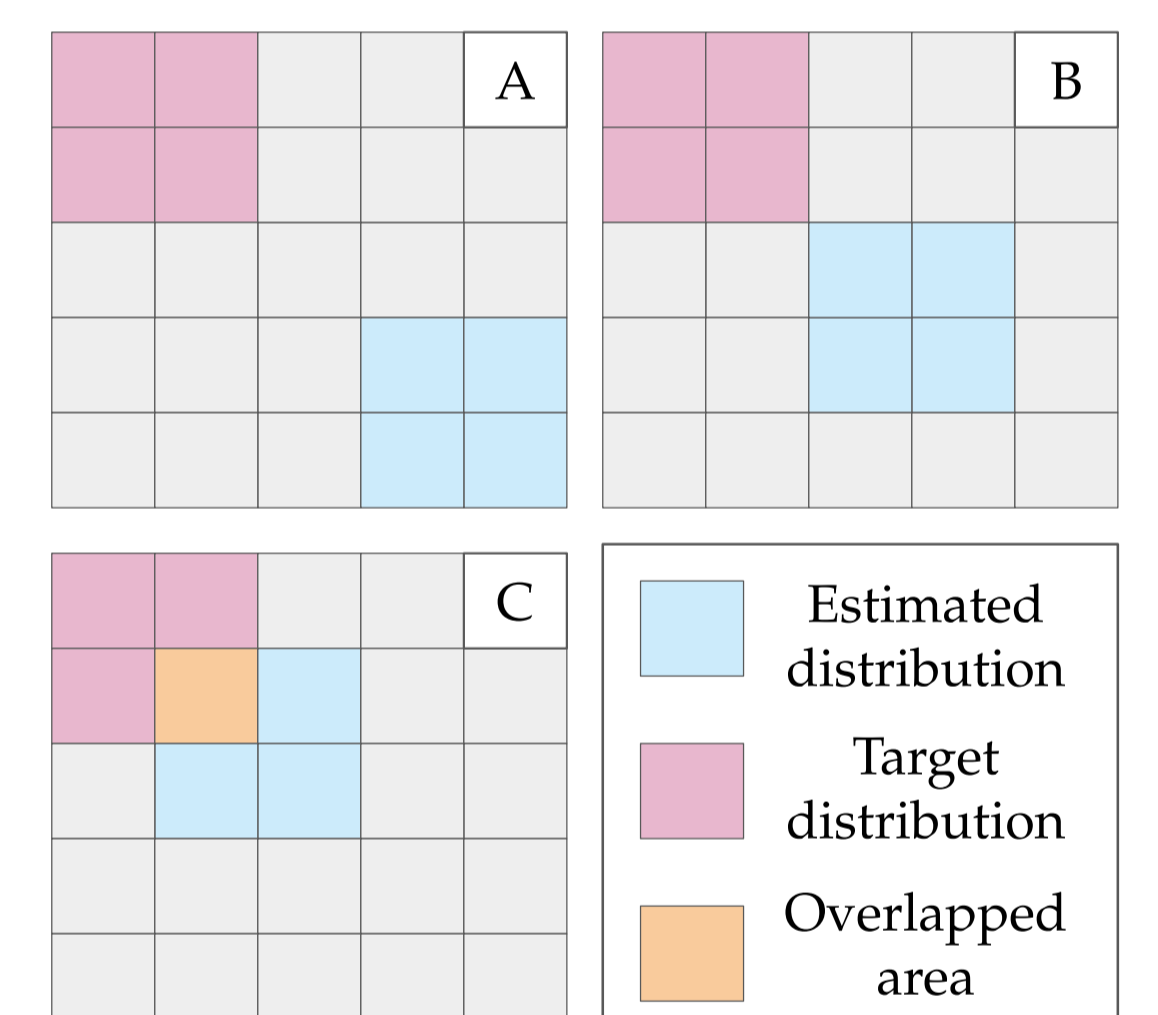
$$H(d, h, w) = \exp\left(-\frac{\sqrt{(d-d_0)^2 + (h-h_0)^2 + (w-w_0)^2}}{2\alpha^2}\right)$$

where  $d, h$  and  $w$  denote a voxel coordinate.  $\alpha$  regulates the distribution spread.

## What should be the optimisation criteria?

Mean Squared Error (MSE) as a cost function works fine at the voxel level, but does not capture global spatial structure. To address this, **Earth Mover's Distance (EMD)** was introduced<sup>6</sup>.

From A to C, the optimal transfer flow would constantly decrease. While, MSE would only decrease in Figure C. This does not account for the similarity of shapes in distributions.



## Heatmap regression suits the task better than coordinate regression to reduce false positive candidates. Essentially, it

- provides probability regions,
- captures correlation and local features,
- shows interpretability on results<sup>5</sup>.

$$\mathcal{E}(p^l, t^l) = -\lambda_1 H_{t^l} \log p^l + \lambda_2 p^l D_{t^l}^\omega$$

Regularised loss    Softmax Cross-Entropy    Distance-weighted spatial enforcer

## Distance matrix as EMD<sup>2</sup> regularisation.

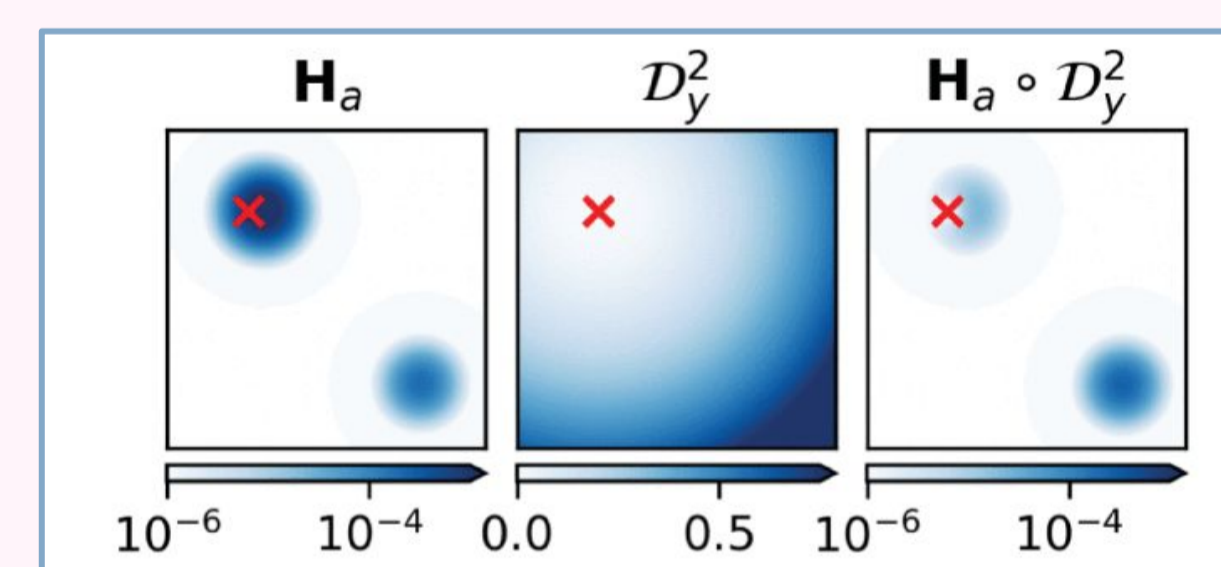


Figure is retrieved directly<sup>6</sup>.

## How to evaluate?

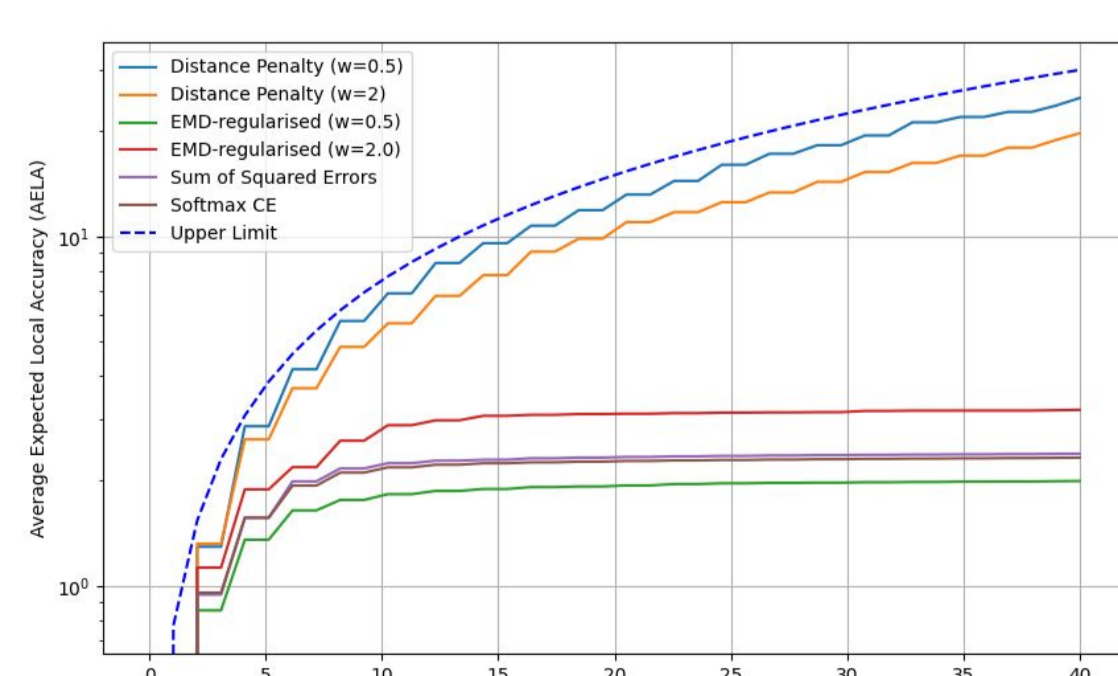
Besides mean distance, we use **Expected Local Accuracy (ELA)** that checks accuracy relative to distance to target.

## We outperform the baseline study!

Our model achieves **lower error scores** than the baseline study<sup>7</sup> where they use Region Proposal Networks (RPNs) formulating "tiny" object detection for **only 5 landmarks**.

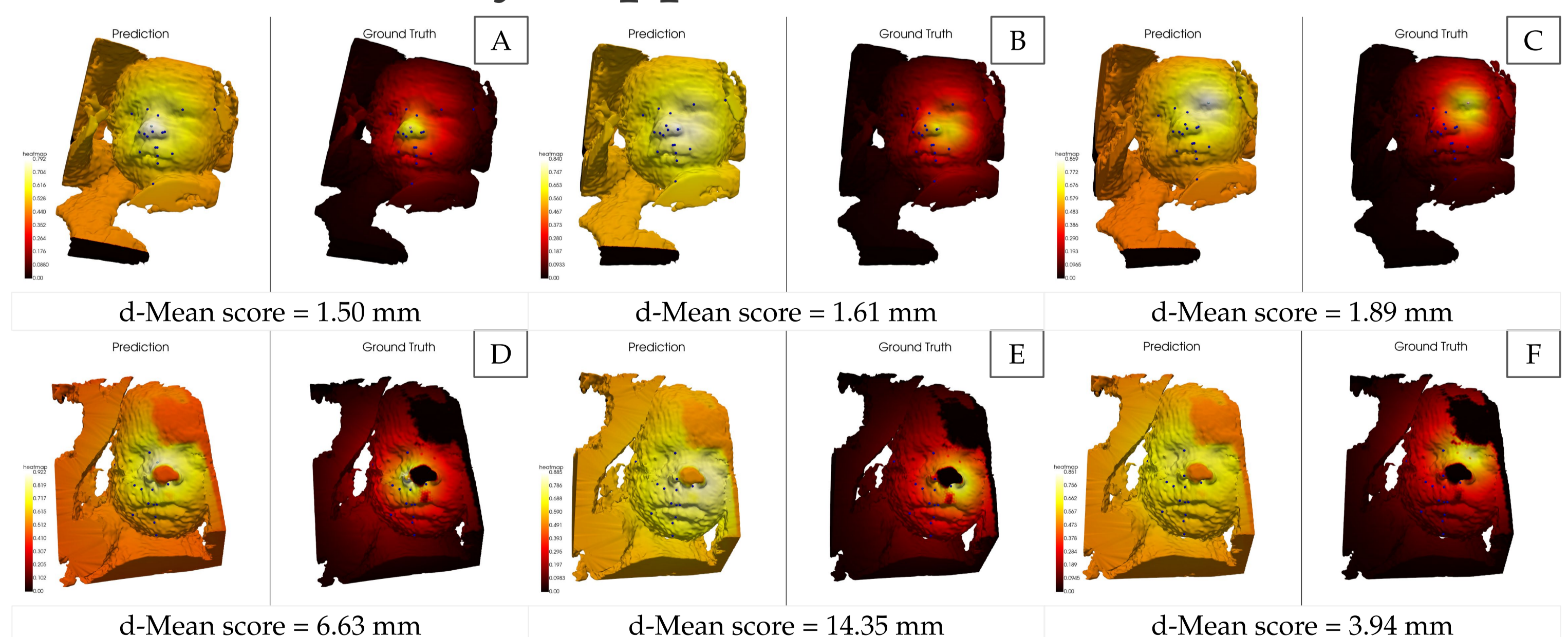
Table of overall distance scores of ResUNet3D trained with loss functions and baseline study. † marks that the setting was subject to cross-fold validation (k=4).

Setting	$w$	d-Mean (mm)
Softmax CE †	-	2.21
Distance Penalty	2	3.24
	0.5	2.37
Sum of Squared Errors	-	2.24
EMD-regularised †	2	3.03
	0.5	1.81
RPN <sup>7</sup>	-	4.38



Expected local accuracy curves for each setting. A random guesser limit was introduced. The chart indicates that within a 40 mm radius, the detector does not change the voxel candidate (argmax) in major misdirection for regularised loss functions.

## What really happens in the "black box"?



A qualitative comparison of model performance on 29-weeks (A-C) and 26-weeks (D-F) fetuses. Target landmarks are prm, aL, and enL (left-to-right). The segmentations are done in 3D Slicer before visualising via PyVista.

## Let's sum up.

Heatmap regression achieves low distance errors and our model can capture facial features while providing clinical interpretation. We have two **limitations**: Data is imbalanced in GA-distribution and the model is prone to error propagation from rotational alignment. **Future work** may focus on enhancing the precision and two-stage learning that includes confidence for landmark visibility.