

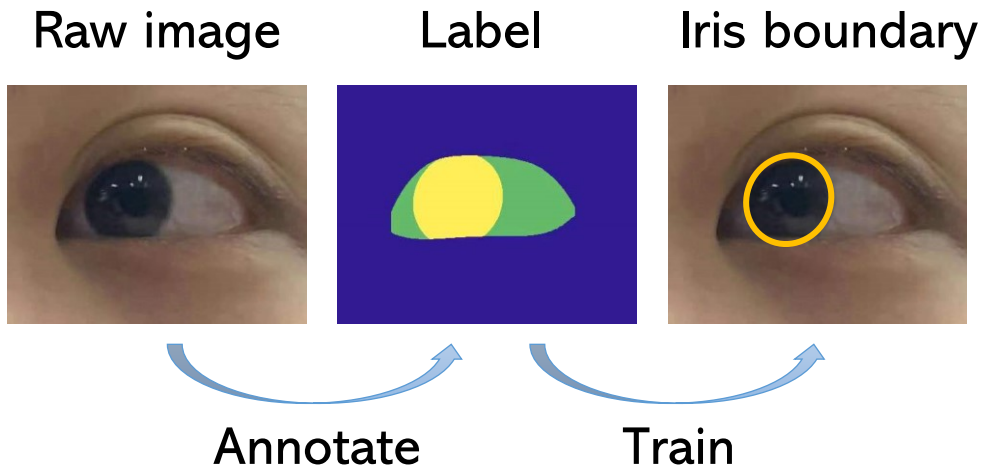
# Practical Gaze Tracking on Any Surface with Your Phone

Jiani Cao<sup>1</sup>, Jiesong Chen<sup>1</sup>, Chengdong Lin<sup>1</sup>, Yang Liu<sup>2</sup>, Kun Wang<sup>1</sup>, Zhenjiang Li<sup>1</sup>  
City University of Hong Kong<sup>1</sup>, University of Cambridge<sup>2</sup>



**Outside the conference version...**

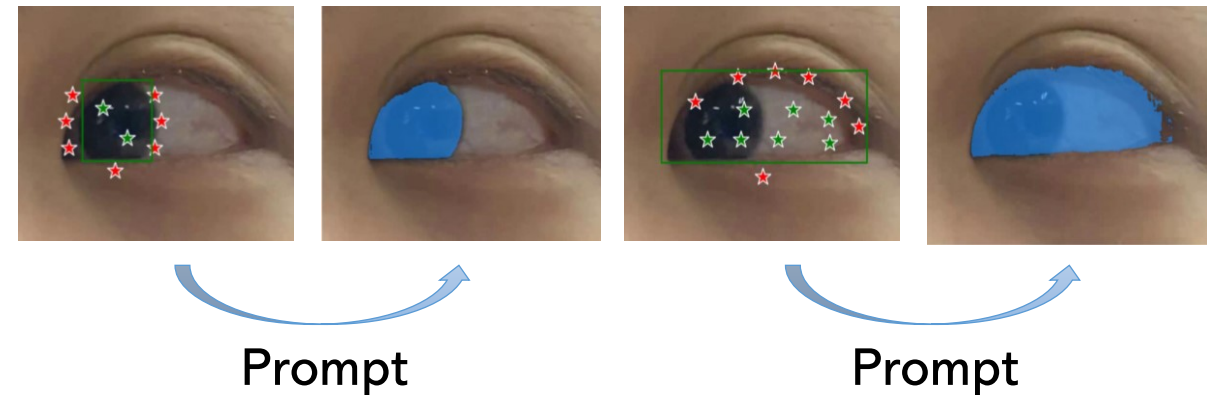
# Training Dataset Label



- Label **quality** directly determines results.
- Manual labelling is **time-consuming**.

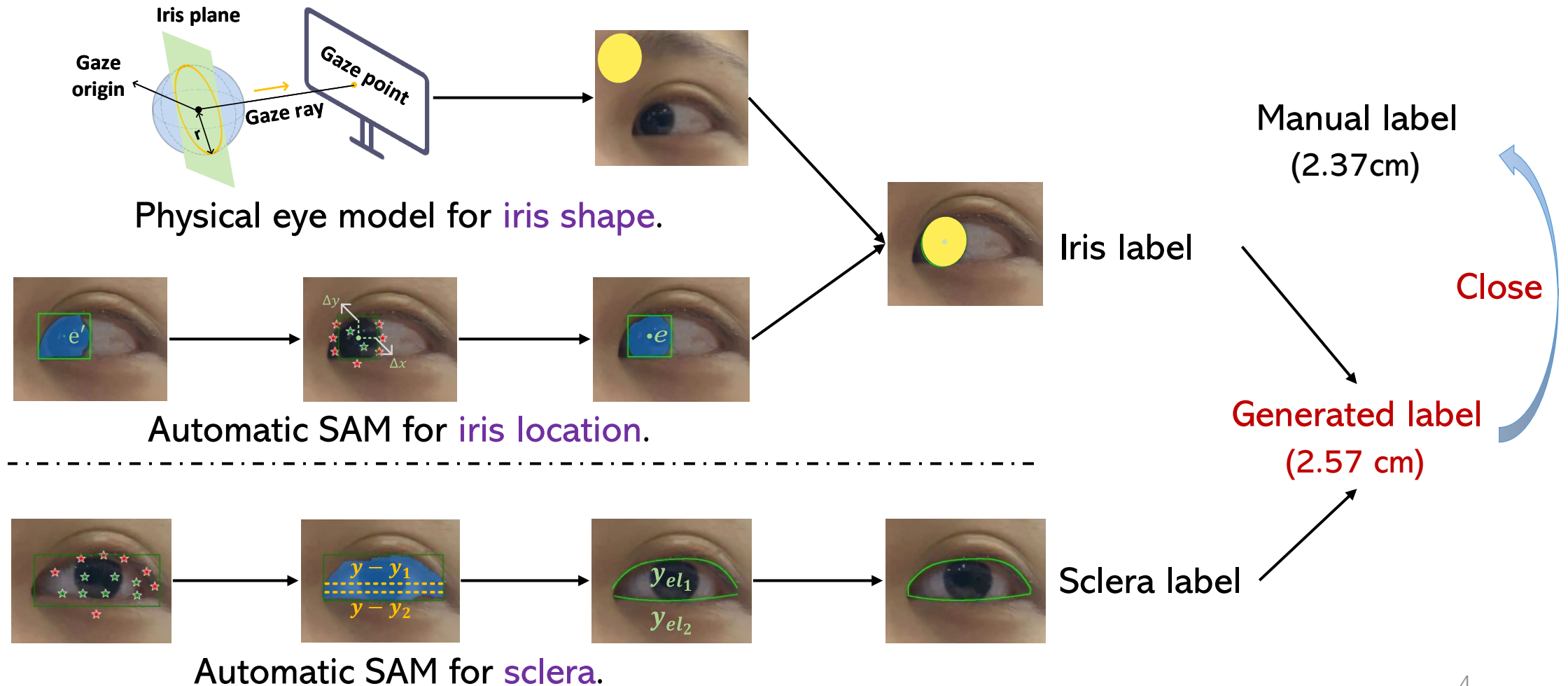
**Foundation model?**  
**Automatic?**

## Using segment anything model (SAM)

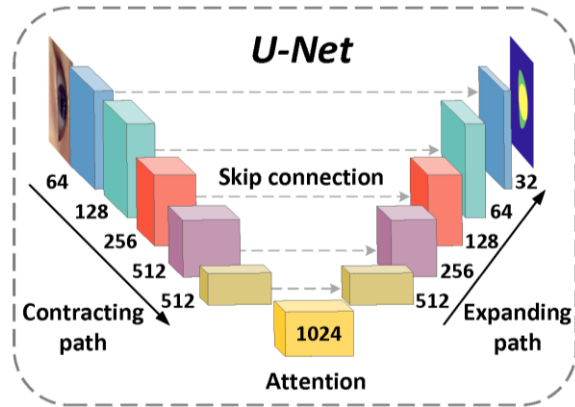


- SAM can generally recognize the iris and sclera.
- However, the edge is quite **unsmooth**. ❌
- The prompt needs to be provided **manually**.

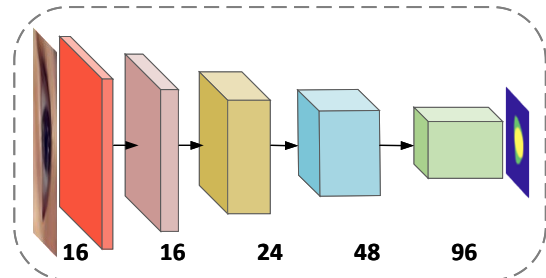
# Training Dataset Label



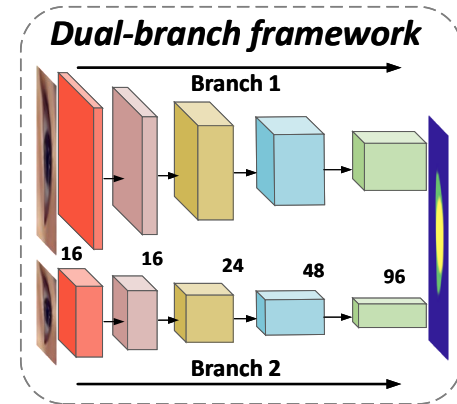
# Latency on Mobile Devices



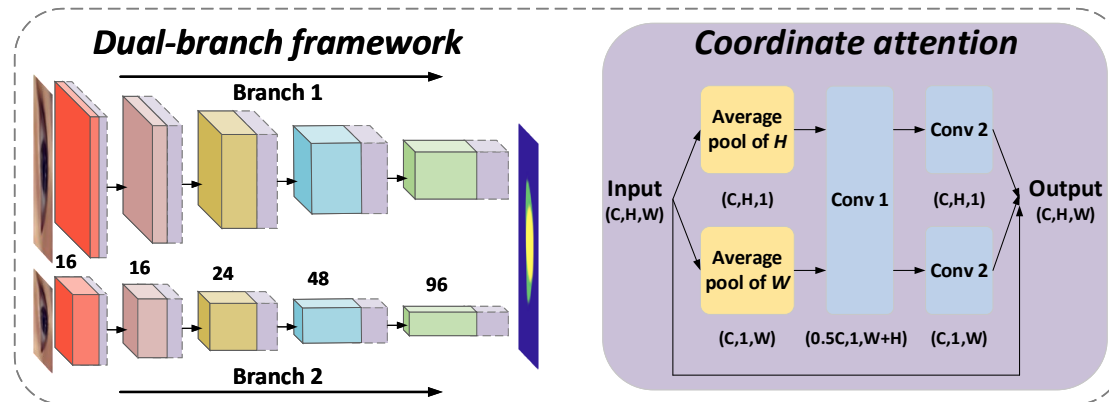
Good performance (2.37cm)  
Bad latency (2.854s)



Bad performance (3.29cm)  
Good latency (0.021s)



Maintain performance (2.97cm)  
Good latency (0.024s)



Further maintain performance (2.87cm)  
Good latency (0.026s)