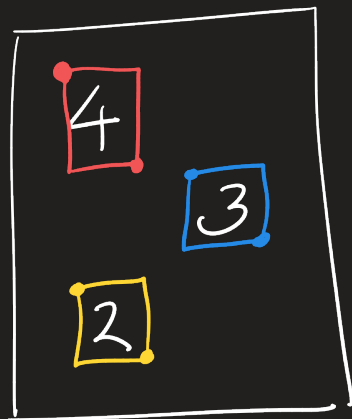


Handwritten Digit object Detection



t_{4top}
 $t_{4bottom}$
 t_{3top}
 $t_{3bottom}$
 t_{2top}
 $t_{2bottom}$
4
3
2

pixel locations for
bounding boxes
for each object.

Labels for all objects
in the image.

Trigger word Classification

Input: Audio file (2-5 seconds)

Record reading
of sentence
Containing trigger
word

"Let's go to the
grocery store"

"stop at traffic
light"

LABEL

go (0)

stop (1)

⋮

