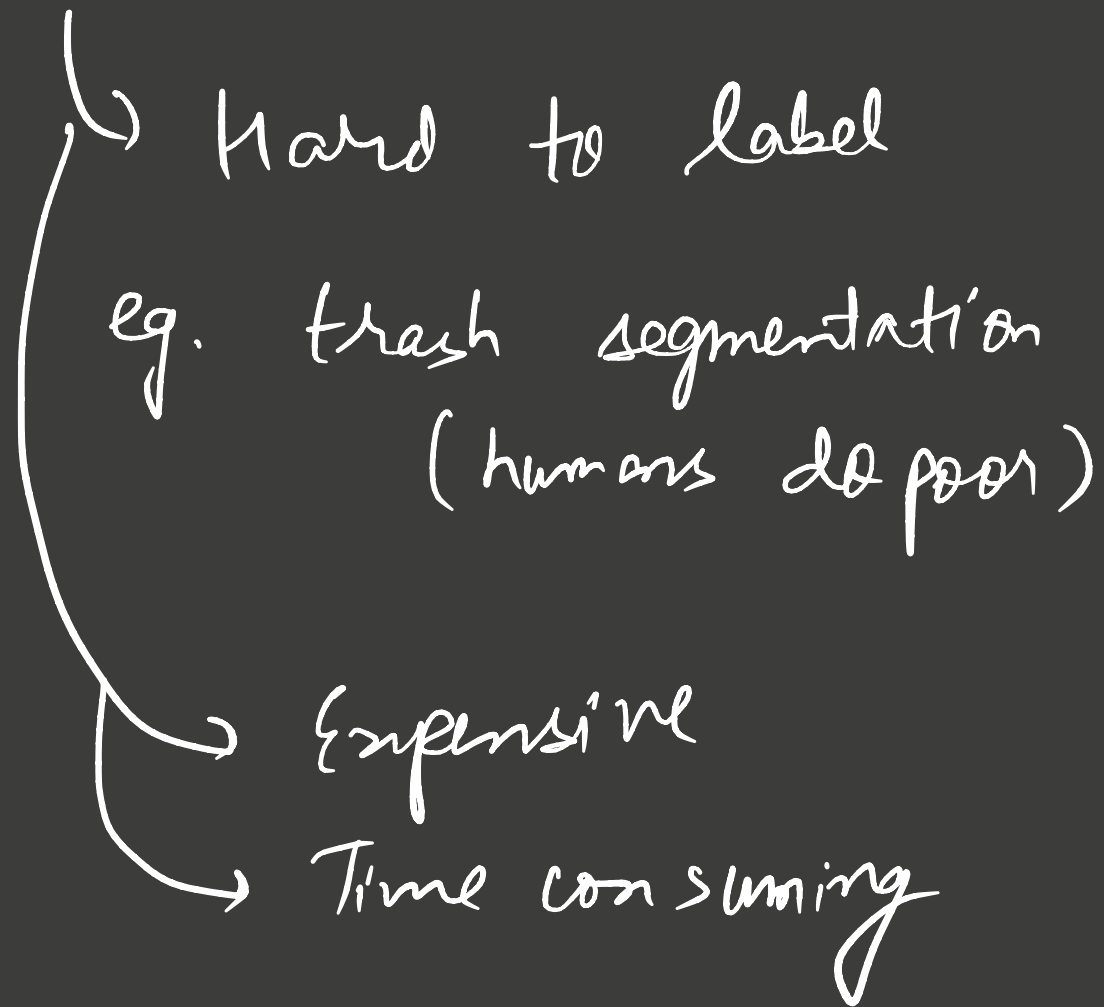


Active learning

Problem with SUPERVISED LEARNING

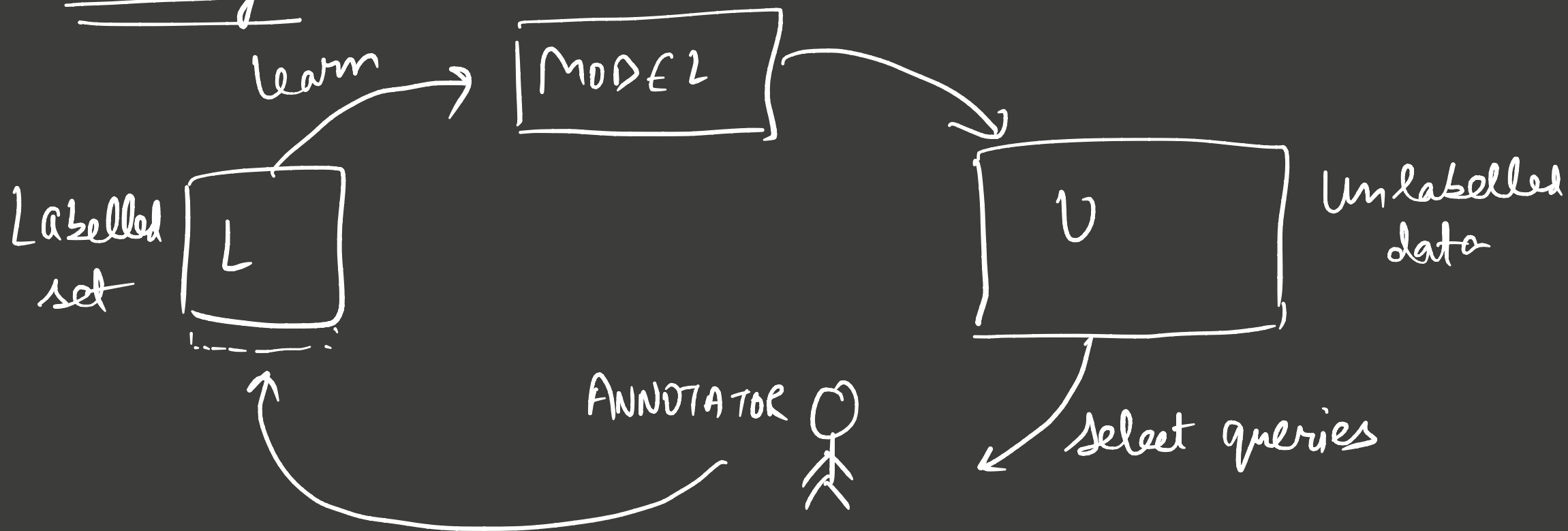
① Require "labelled" data



- ① Unlabeled data is easy to get / cheap
- ② Label 'as few' ~~or~~ points as possible.

Budget
Accuracy desired

Setting

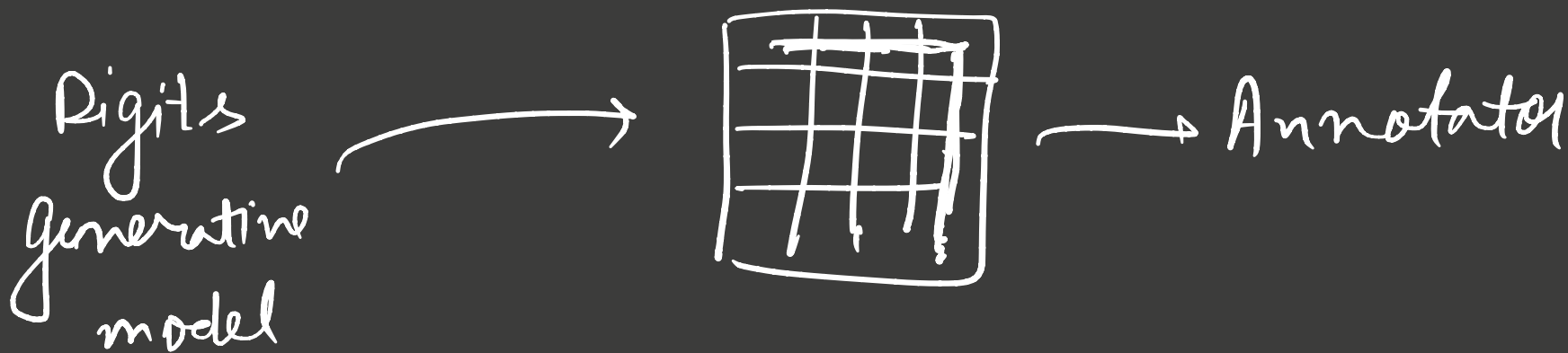
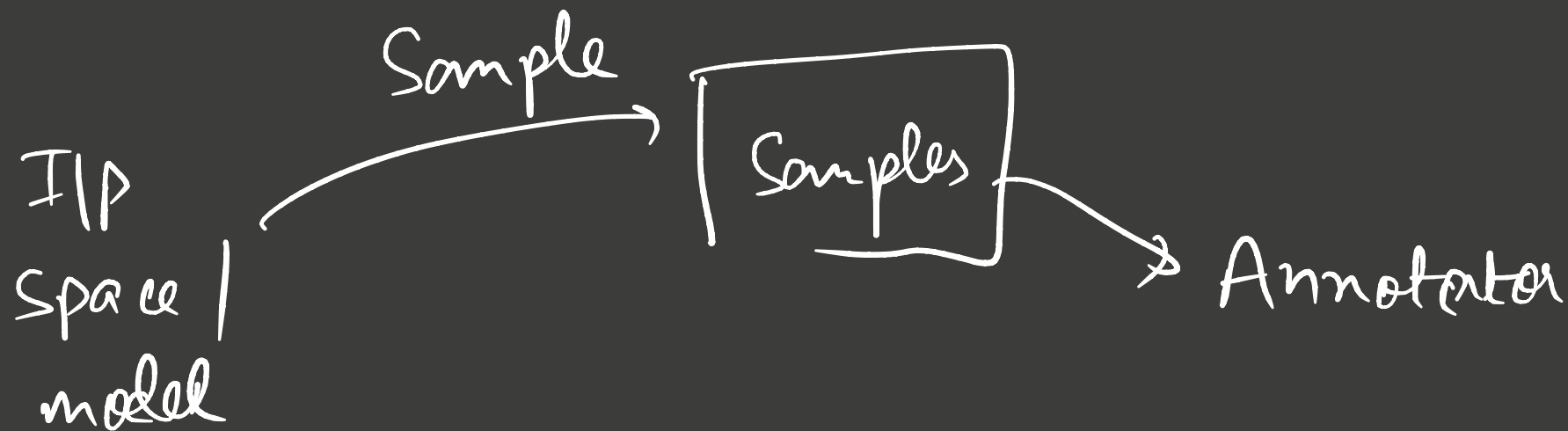


Evaluation of Active Learning



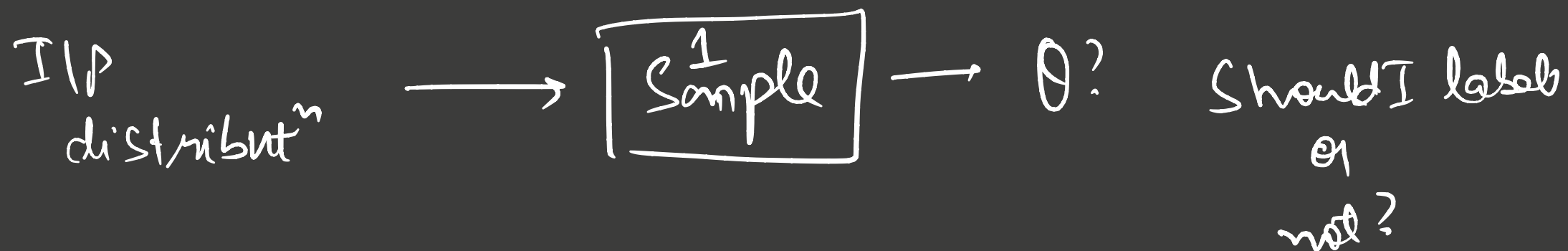
Scenarios of Active Learning

① Membership Query Synthesis

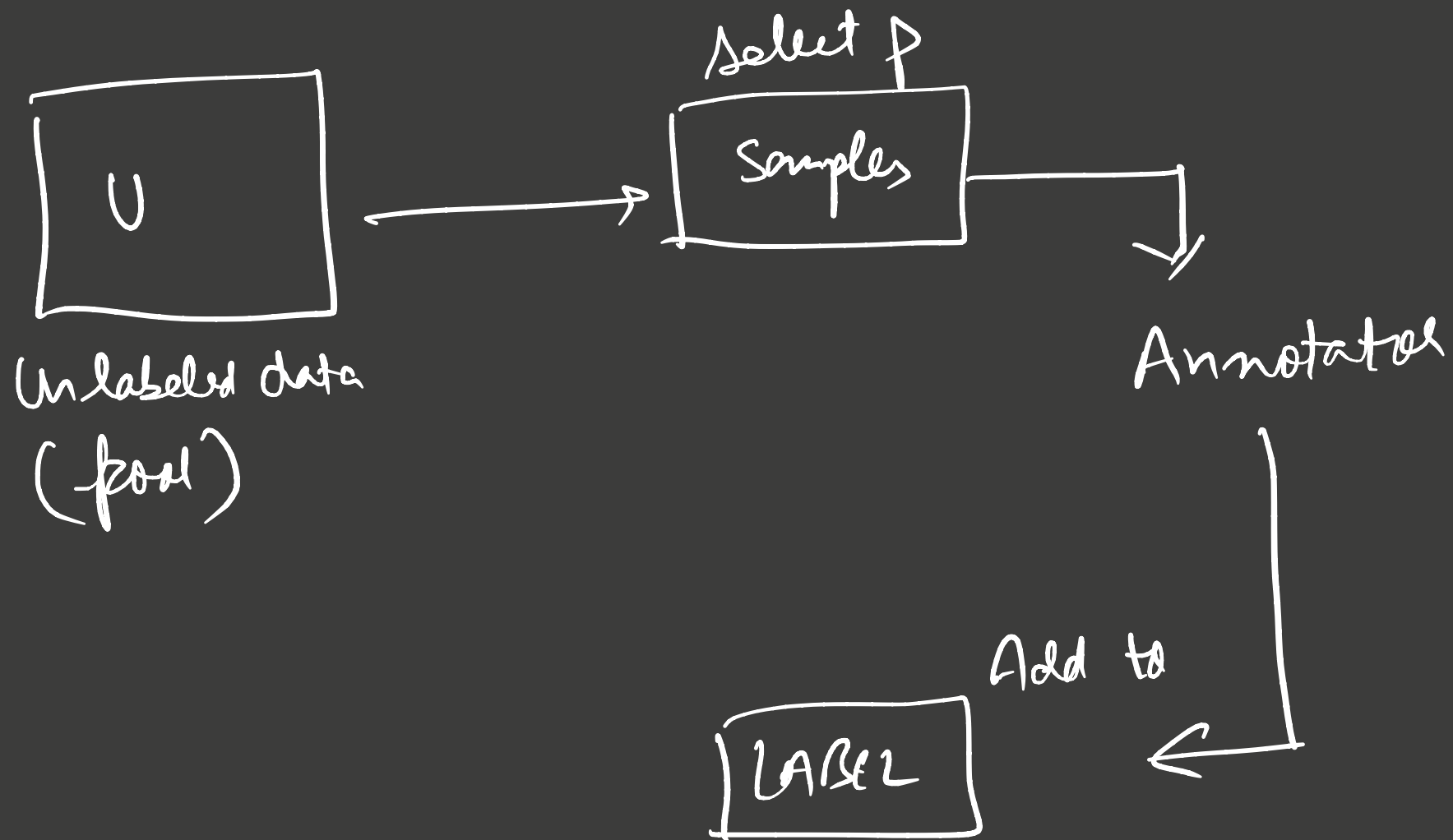


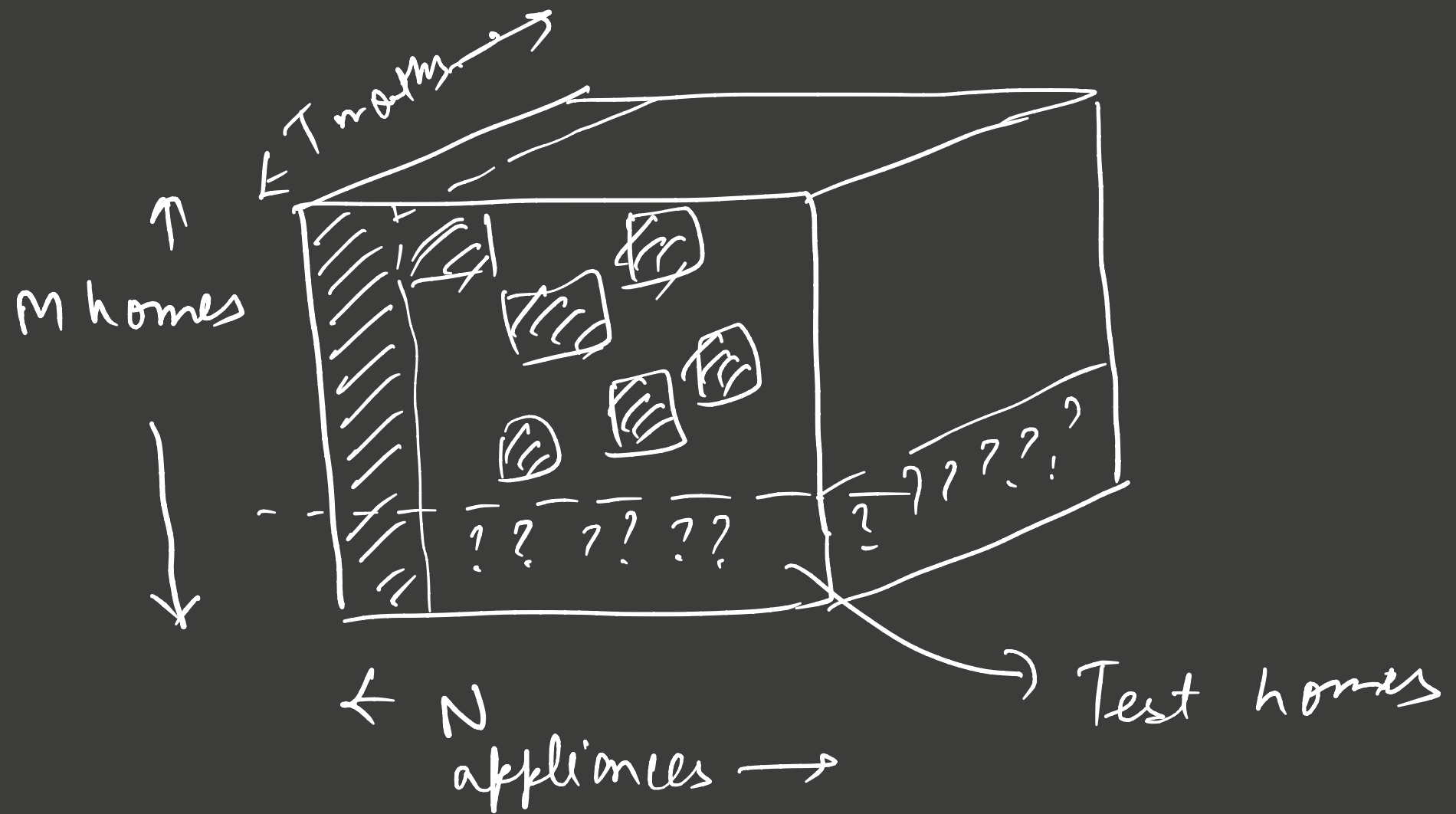
Can lead to incoherent examples.

② Stream based



③ Pool based.





Q) Given budget of X sensors
 Choose $\langle \text{home, appliance} \rangle \Rightarrow$
 Maximize accuracy.

Query Strategies (for Pool based Sampling)

(1) Uncertainty sampling

→ Query points you're most uncertain about

→ least confident estimate

$$x_{LC}^* = \operatorname{argmax}_x (1 - P_{\theta}(\hat{y}|x))$$

\hat{y} = class with max. probability

→ Margin sampling

$$x_M^* = \underset{x}{\operatorname{argmin}} (P_{\theta}(\hat{y}_1 | x) - P_{\theta}(\hat{y}_2 | x))$$

\hat{y}_1, \hat{y}_2 are the two
most probable
classes

Q) How will you do uncertainty sampling for KNN
(KNN
(classifier))

Votes (+) ~ Votes (-)

Query by committee (QBC)

1) Create a committee of models $C = \{\theta^1, \theta^2, \dots, \theta^c\}$

trained on labeled set (L) but represent

Competing hypothesis,

2) Each of the model votes

3) Most informative instance = most disagreement
in votes
(entropy high)

eg. θ^{BC} with Random Forest

$\theta^1 =$ Model 1 = R.F. with seed 1

$\theta^2 =$ " 2 = " " " 2

⋮

eg. θ^{BC} with KNN

$\theta^1 =$ KNN with $K=1$

$\theta^2 =$ - - - - -

OBC for regression

⇒ Compute variance amongst committee

⇒ choose instances with highest variance.