**Query**

What did I put in the black bin in the kitchen?

VLM (8 frames)

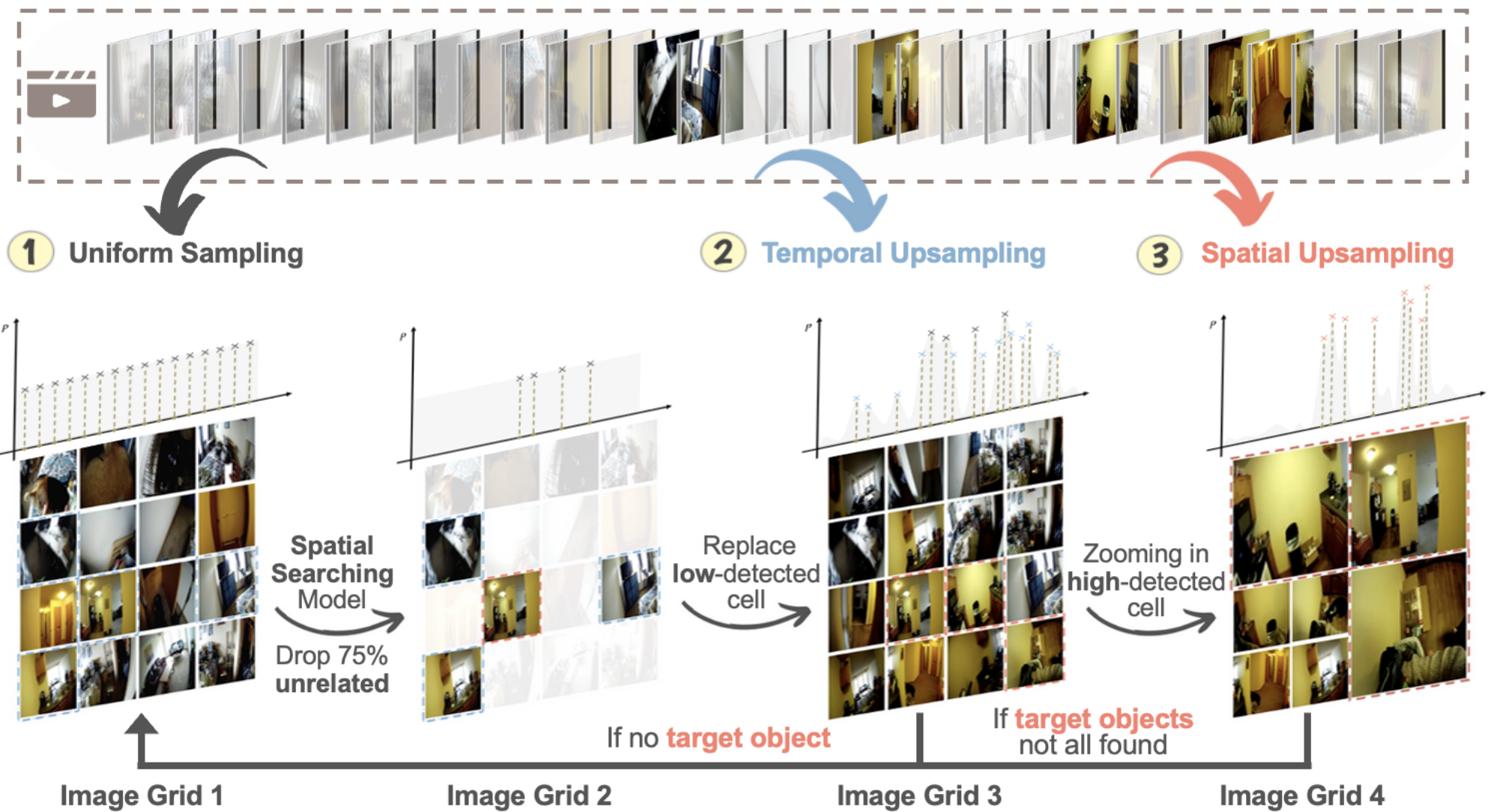**Target Object**
The black bin in the kitchen

**Cue Object**
Refrigerator, Toaster, Microwave, Tables, Stovetop, Stools, Floors, ...

**Grounding**

① Uniform Sampling
② Temporal Upsampling
③ Spatial Upsampling

**Spatial Searching Model**
Drop 75% unrelated

Replace **low**-detected cell

Zooming in **high**-detected cell

Image Grid 1    Image Grid 2    Image Grid 3    Image Grid 4

If no **target object**

If **target objects** not all found
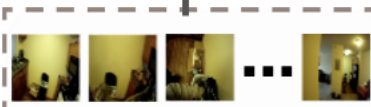
**Iterative Temporal Searching**

**Answer**

"vacuum dust"

VLM (K frames)

If confirmed **target objects**

Query

Selected Frames

**Downstream QA**