**Input:** Multimodal user-imagined *IDEA* to generate

*IDEA* 1:

photo of **Bill Gates** with the **same hand gesture** as in the given image
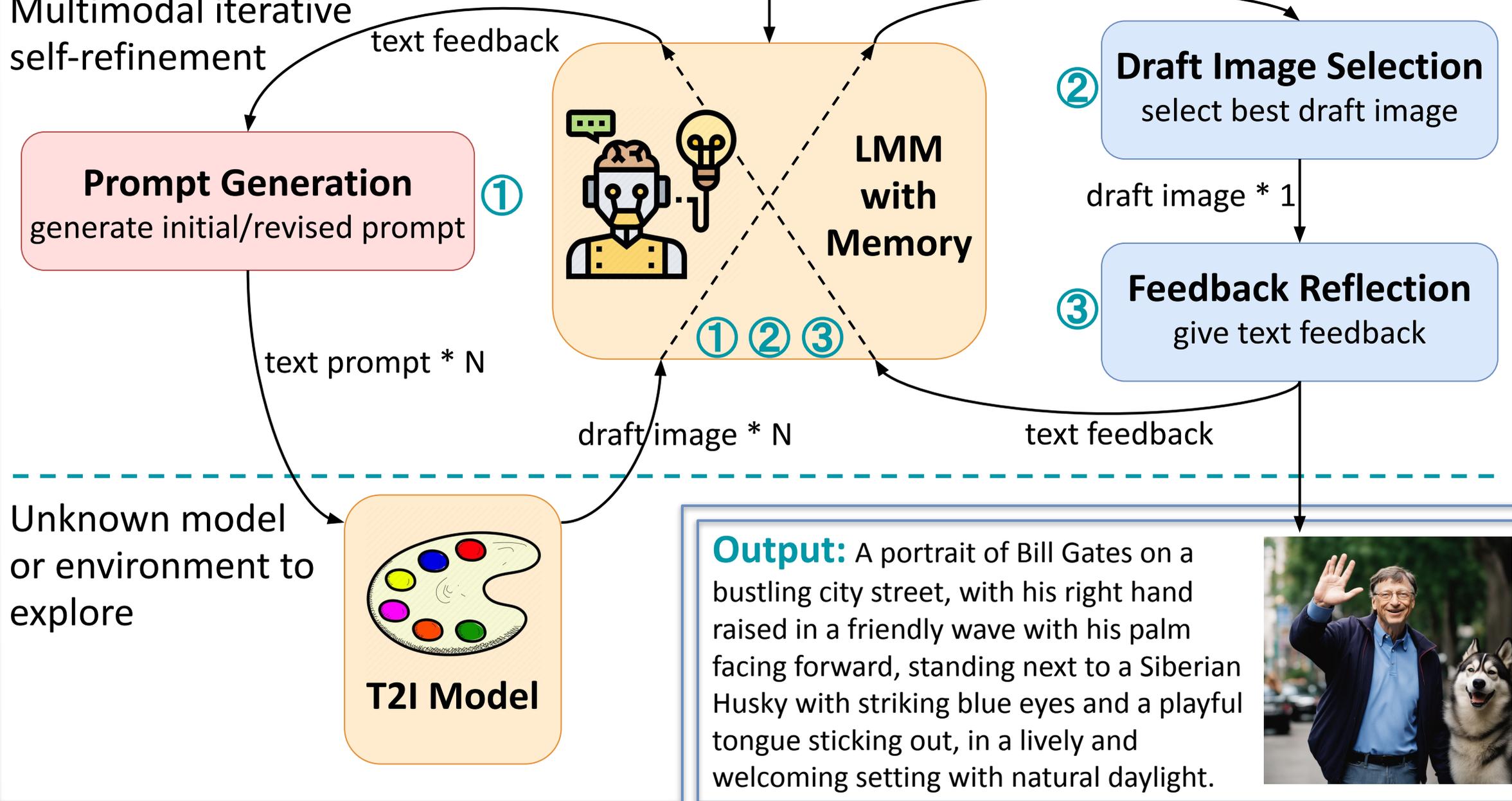


, with a dog looks like this one in the image



*Idea2Img Framework*

Multimodal iterative self-refinement

*IDEA*

draft image * N

text feedback

**Prompt Generation**
generate initial/revised prompt ①

LMM with Memory

① ② ③

② **Draft Image Selection**
select best draft image

text prompt * N

draft image * 1

③ **Feedback Reflection**
give text feedback

Unknown model or environment to explore

draft image * N

text feedback

**T2I Model**

**Output:** A portrait of Bill Gates on a bustling city street, with his right hand raised in a friendly wave with his palm facing forward, standing next to a Siberian Husky with striking blue eyes and a playful tongue sticking out, in a lively and welcoming setting with natural daylight.



*IDEA* 2:

photo of **Bill Gates** with the **same suit** as in the given image on the street

, with a dog looks like this one in the image

*Idea2Img*

*IDEA* 3:

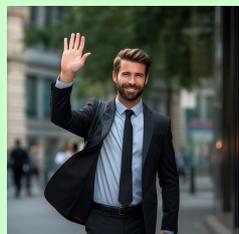**cartoon** drawing of the person as in the given image playing with a dog on the beach

, with a dog looks like this one in the image

*Idea2Img*