**What is this a picture of ?**
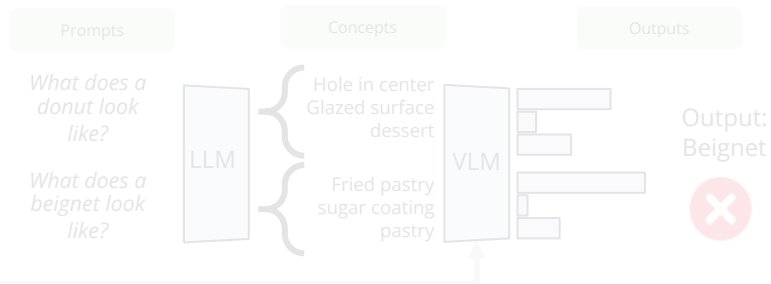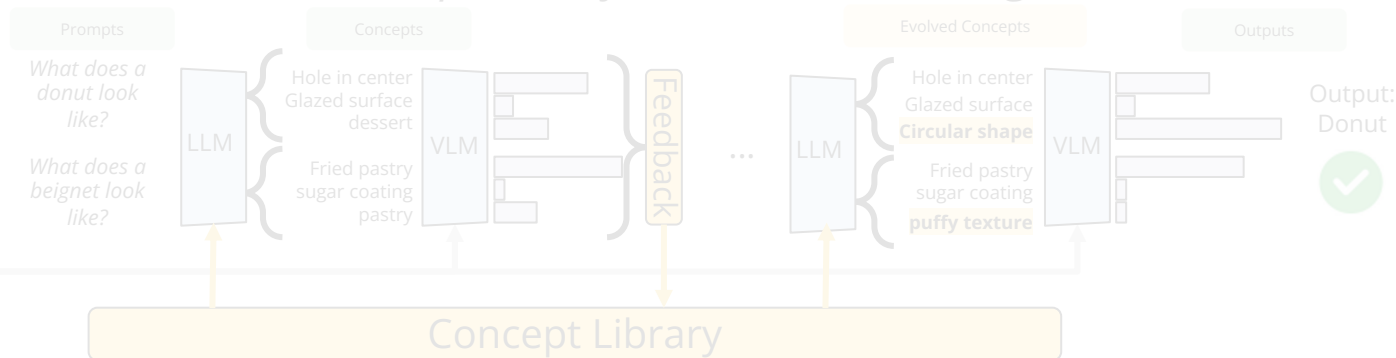
A donut (🍩) or a beignet (🥮)?

**Prior work**: LLMs generate *visual features* in **one shot**.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

LLM

Concepts

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Outputs

Output: Beignet ❌

*However, the LLM's proposed features is different from the VLM's sensitivity to features.*

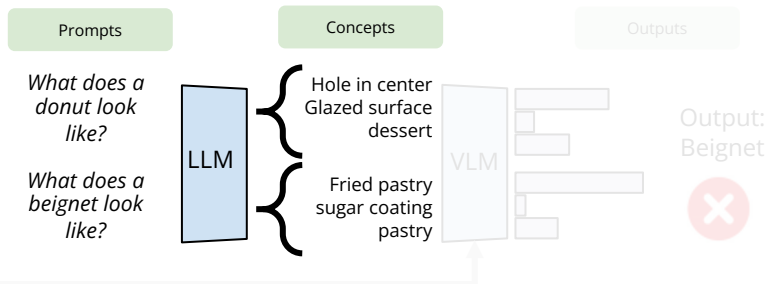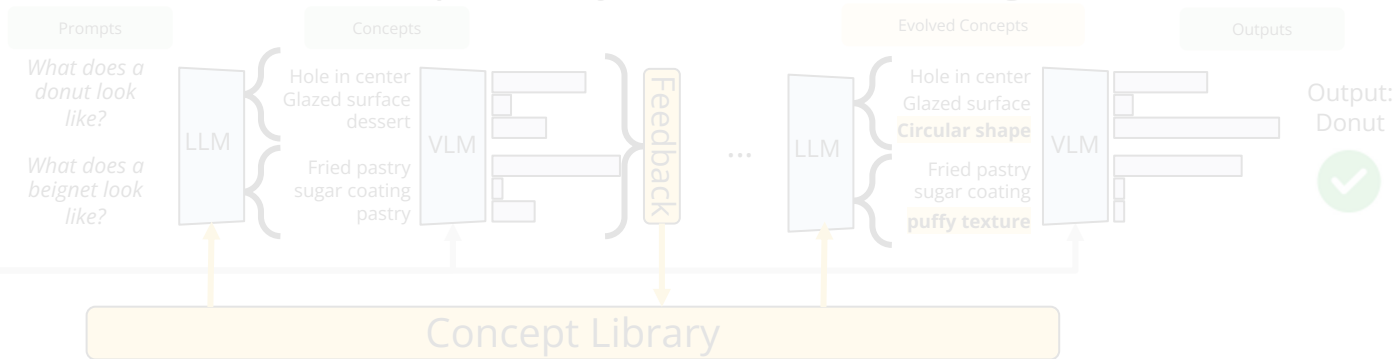Our work: a *concept library* that **evolves** using the VLM as a critic.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

LLM

Concepts

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Feedback

...

LLM

Evolved Concepts

Hole in center
Glazed surface
**Circular shape**

Fried pastry
sugar coating
**puffy texture**

VLM

Outputs

Output: Donut ✅

Concept Library

**Prior work**: LLMs generate *visual features* in **one shot**.

What is this a picture of ?

A donut (🍩) or a beignet (🍮)?

Prompts

*What does a donut look like?*

*What does a beignet look like?*

LLM

Concepts

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Outputs

Output: Beignet

❌

However, the LLM's proposed features is different from the VLM's sensitivity to features.

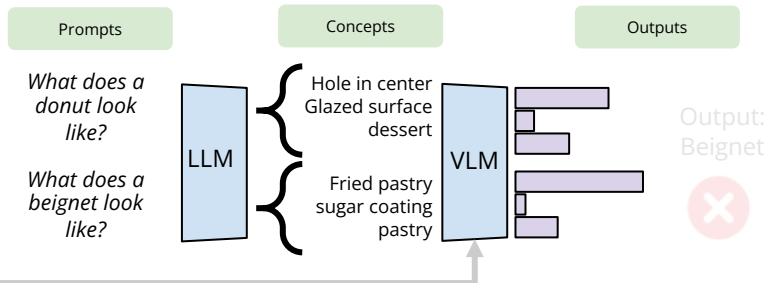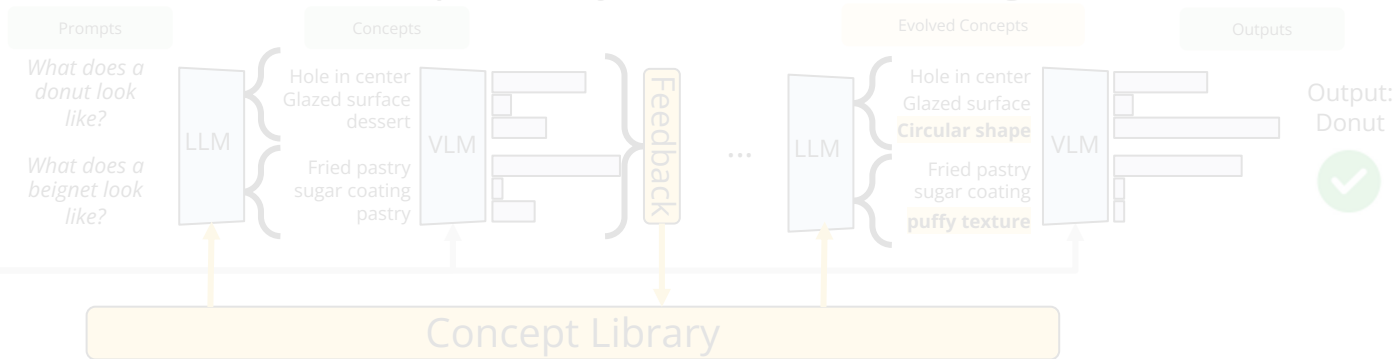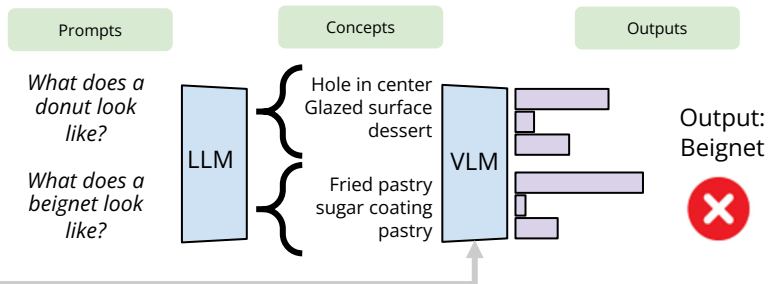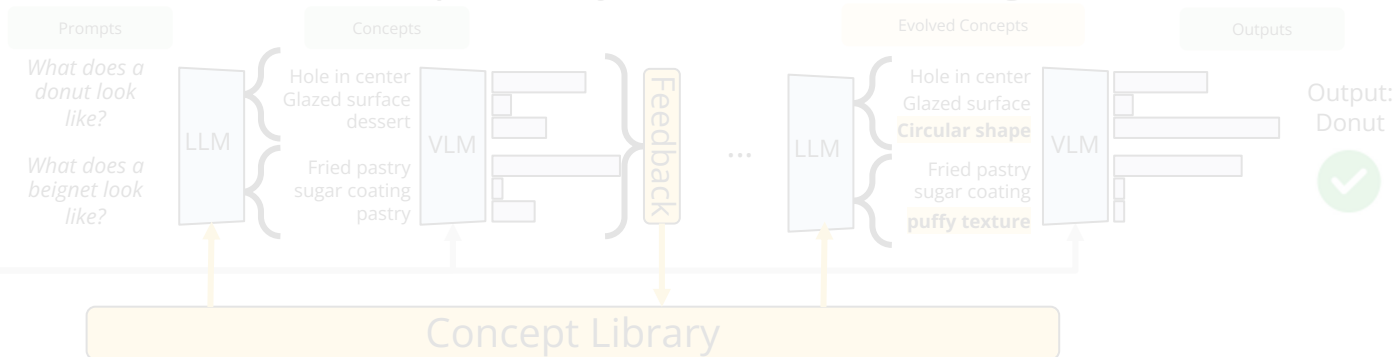Our work: a *concept library* that **evolves** using the VLM as a critic.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

LLM

Concepts

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Feedback

...

LLM

Evolved Concepts

Hole in center
Glazed surface
**Circular shape**

Fried pastry
sugar coating
**puffy texture**

VLM

Outputs

Output: Donut

✓

Concept Library

**What is this a picture of ?**

A donut (🍩) or a beignet (🥮)?

Prior work: LLMs generate *visual features* in **one shot**.

| Prompts | Concepts | Outputs |

*What does a donut look like?*

LLM

Hole in center
Glazed surface
dessert

VLM

Output: Beignet ❌

*What does a beignet look like?*

Fried pastry
sugar coating
pastry

*However, the LLM's proposed features is different from the VLM's sensitivity to features.*

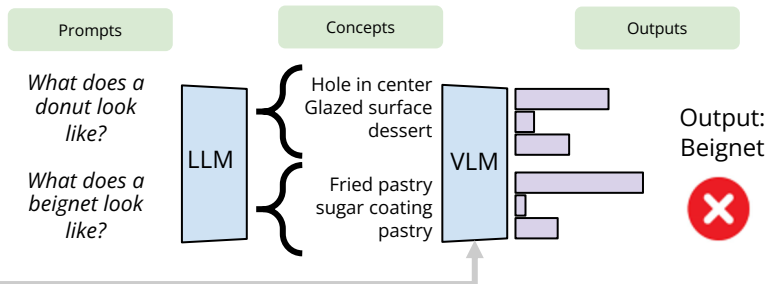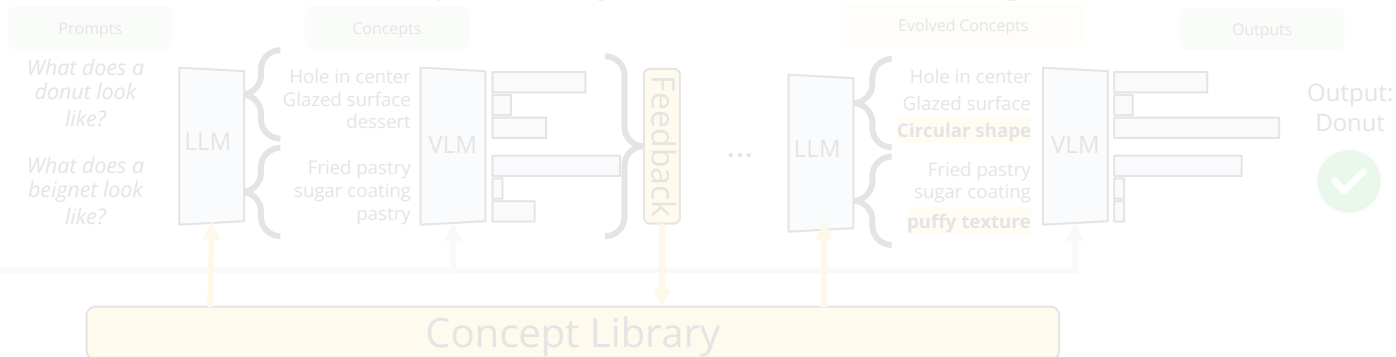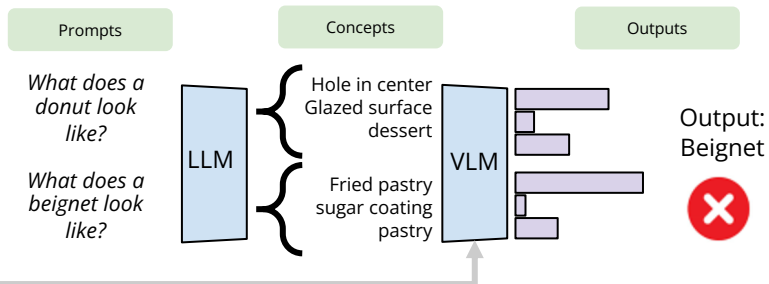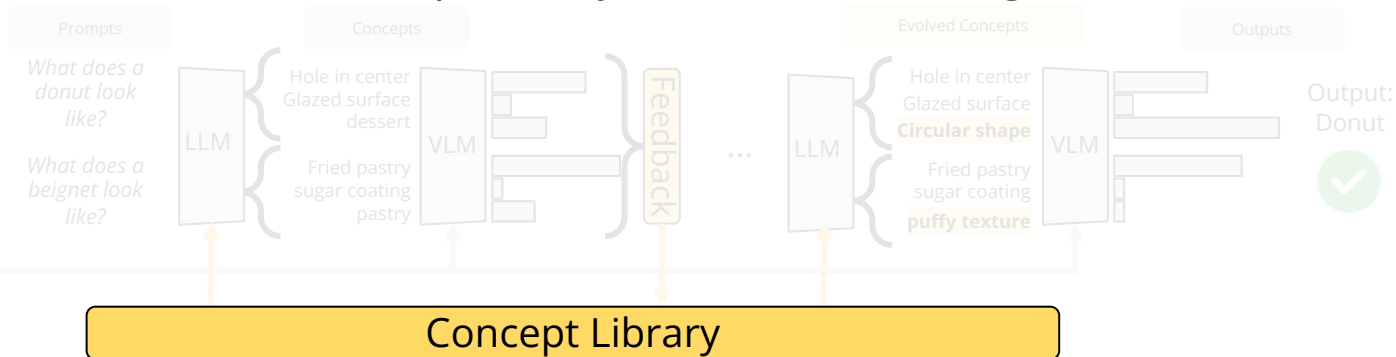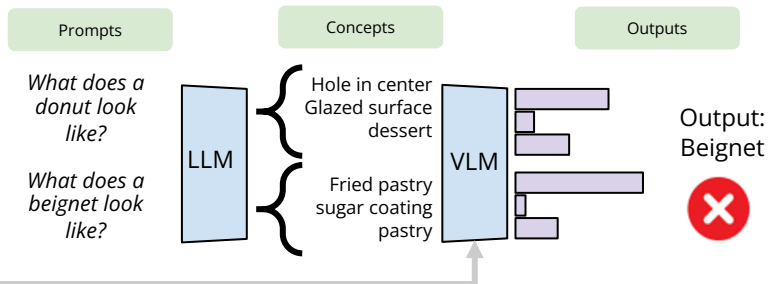Our work: a *concept library* that **evolves** using the VLM as a critic.

| Prompts | Concepts | Evolved Concepts | Outputs |

*What does a donut look like?*

LLM

Hole in center
Glazed surface
dessert

VLM

Feedback

... LLM

Hole in center
Glazed surface
**Circular shape**

VLM

Output: Donut ✓

*What does a beignet look like?*

Fried pastry
sugar coating
pastry

Fried pastry
sugar coating
**puffy texture**

Concept Library

**What is this a picture of ?**

A donut (🍩) or a beignet (🥯)?

Prior work: LLMs generate *visual features* in **one shot**.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Outputs

Output:
Beignet

❌

*However, the LLM's proposed features is different from the VLM's sensitivity to features.*

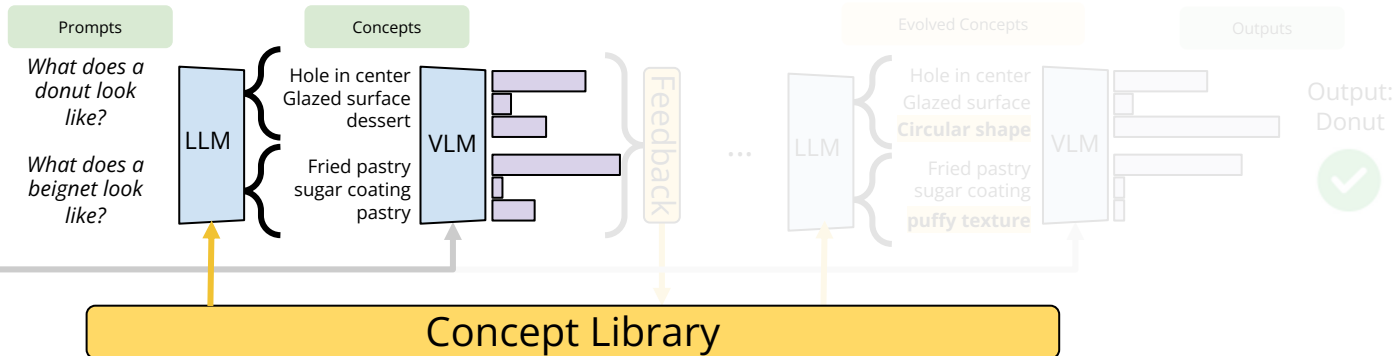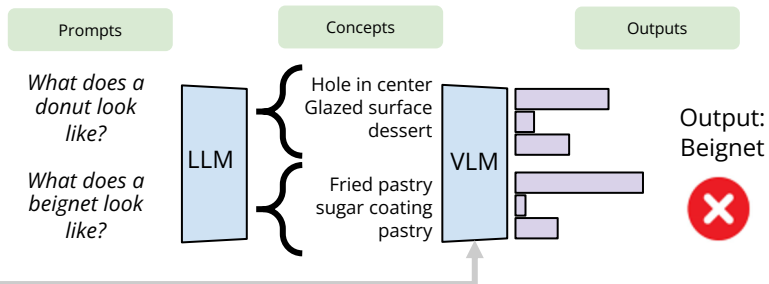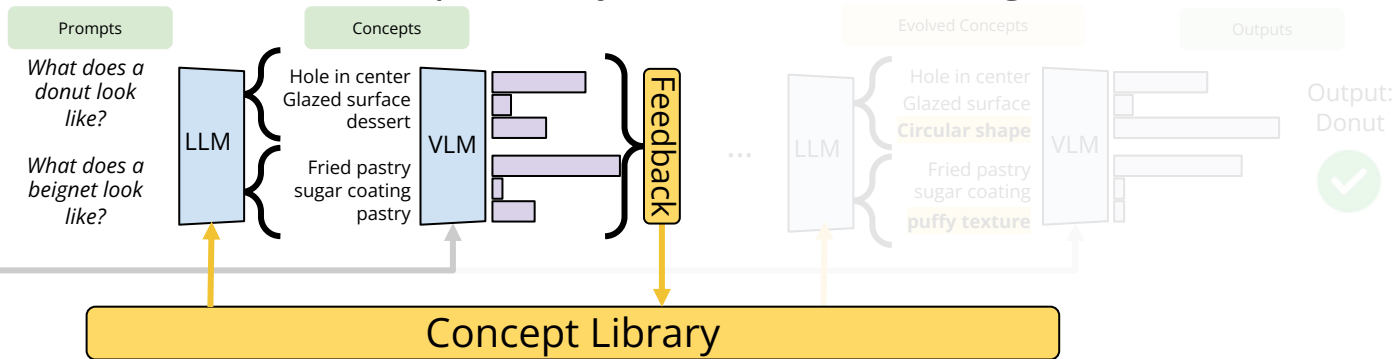Our work: a *concept library* that **evolves** using the VLM as a critic.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Feedback

...

Evolved Concepts

LLM

Hole in center
Glazed surface
**Circular shape**

Fried pastry
sugar coating
**puffy texture**

VLM

Outputs

Output:
Donut

✅

Concept Library

*Pseudo* Confusion Matrix

*Pseudo* Confusion Matrix

Please suggest visual features that distinguish **a slaty backed gull** from **a california gull**.

The current descriptors for **slaty backed gull** are: {cls1_concepts}.
The current descriptors for **california gull** are: {cls2_concepts}.

$\downarrow$

**LLM**

$\downarrow$

```json
{
  "reasoning": [
          "Slaty-backed gulls have a darker, slate-gray back and wings, whereas California gulls have a lighter gray back and wings, making the back coloration a key distinguishing feature.",
          "Slaty-backed gulls exhibit prominent white 'mirrors' on their primary feathers, which are larger and more distinct compared to those of the California gull.", ...
          ],
  "features": [
          "Slate-gray back and wings",
          "Large white mirrors on primary feathers",
...
          ]
}
```



*Pseudo* **Confusion Matrix**

**What is this a picture of ?**

A donut (🍩) or a beignet (🥯)?

Prior work: LLMs generate *visual features* in **one shot**.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Outputs

Output:
Beignet

❌

*However, the LLM's proposed features is different from the VLM's sensitivity to features.*

Our work: a *concept library* that **evolves** using the VLM as a critic.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Feedback

...

Evolved Concepts

LLM

Hole in center
Glazed surface
**Circular shape**

Fried pastry
sugar coating
**puffy texture**

VLM

Outputs

Output:
Donut

✅

Concept Library
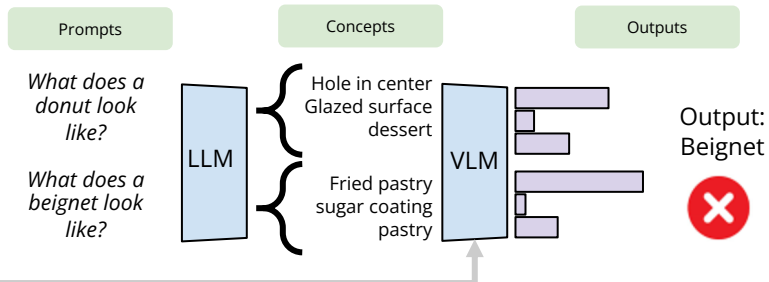
**What is this a picture of ?**

A donut (🍩) or a beignet (🍤)?

Prior work: LLMs generate *visual features* in **one shot**.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Outputs

Output:
Beignet

❌

*However, the LLM's proposed features is different from the VLM's sensitivity to features.*

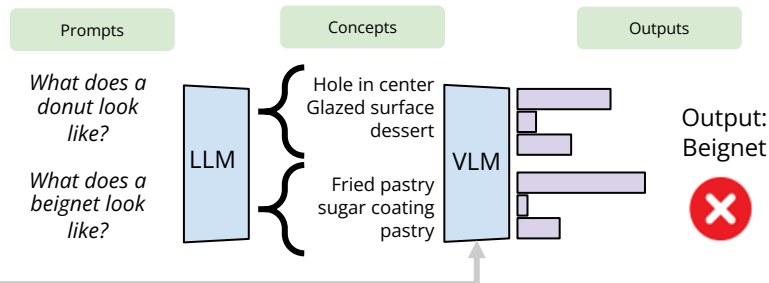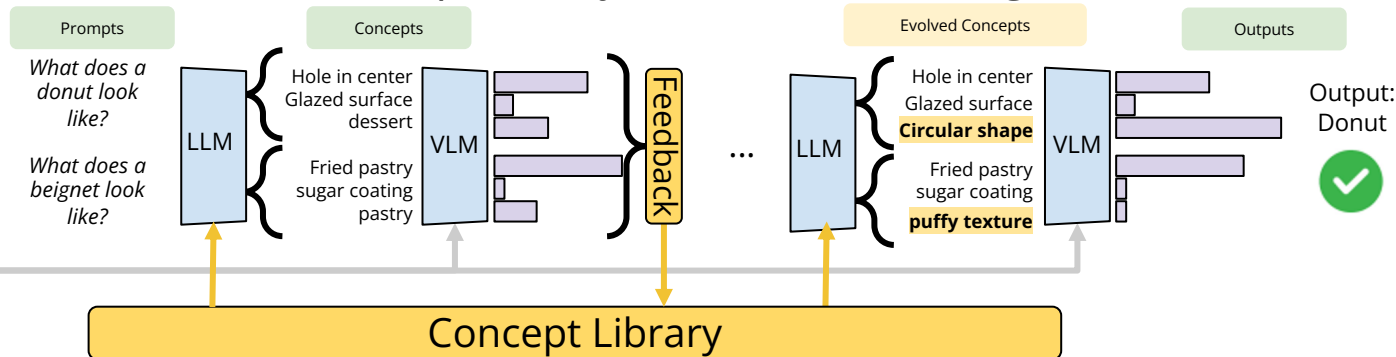Our work: a *concept library* that **evolves** using the VLM as a critic.

Prompts

*What does a donut look like?*

*What does a beignet look like?*

Concepts

LLM

Hole in center
Glazed surface
dessert

Fried pastry
sugar coating
pastry

VLM

Feedback

...

Evolved Concepts

LLM

Hole in center
Glazed surface
**Circular shape**

Fried pastry
sugar coating
**puffy texture**

VLM

Outputs

Output:
Donut

✅

Concept Library

**Key insight:** The concept library is simultaneously bootstrapping the VLM inputs and the LLM inputs.