**(a) Rendering Method**

**CoT Plain Text**

Weng earns 12/60 = $<<12/60=0.2>>0.2 per minute. Working 50 minutes, she earned 0.2 x 50 = $<<0.2*50=10>>10.

**Rendering** — Font Size | Text Color | Image Height | …

**Rendered Image**

Weng earns 12/60 = $<<12/60=0.2>>0.2 per minute. Working 50 minutes, she earned 0.2 x 50 = $<<0.2*50=10>>10.

**Vision Encoder** ❄

**Question Text**

Weng earns $12 an hour for babysitting. Yesterday, she just did 50 minutes of babysitting. How much did she earn?

<vision embedding>

$\mathcal{L}_{MSE}$

**10** Answer

$\mathcal{L}_{CE}$

<answer>

<question>

**Visual Projection Head**

**LLM Backbone**

<|img_begin|>    <latent reasoning>    <|img_end|>

**(b) Latent Reasoning Method**