# A  QUALITATIVE ANALYSIS: DISTILLING SUCCESS AND FAILURE

To validate the robustness of our thought compression, we analyze both success and failure scenarios. Our analysis reveals that WISE-S achieves high efficiency without introducing new reasoning errors, maintaining high **Semantic Fidelity** to the original verbose reasoning.

**Success Cases: Concise Rationale as a Sufficient Summary.** As shown in Figure 1, in successful instances, WISE-S demonstrates the ability to distill the *decision logic* while discarding *visual redundancy*.

- *Attribute Grounding:* In the equestrian example (Figure 1, Bottom-Right), the detailed explanation ($\tau_d$) engages in a verbose verification process, checking the horse's position and the nature of the sport. In contrast, WISE-S ($\tau_c$) directly extracts the discriminative features—"red and white obstacle" and "foreground"—which are sufficient to localize the mask.

- *Functional Reasoning:* In the cave exploration example (Figure 1, Bottom-Left), the model must process a negative constraint ("did not consider diving"). WISE-S correctly reasons that the target area must be "above the water level," efficiently pruning the search space without the need for the extensive geological description found in the detailed chain.

These examples confirm our hypothesis that the concise rationale learns to act as a sufficient summary for the final answer, effectively bridging the gap between instruction and segmentation.

**Failure Cases: Consistent Limitations.** We closely examined cases with low IoU (Figure 2) to determine if the brevity constraint caused the failure. Interestingly, we found that **the compression mechanism is rarely the culprit**.

- **Fidelity in Failure:** In the "Warthog" example, the concise rationale ($\tau_c$) correctly identifies the target object as "tusks," fully capturing the semantic core of the verbose explanation ($\tau_d$). The failure to segment the specific tusks (likely masking the whole face) is a shared limitation in *spatial grounding* inherent to the base Vision-Language Model, occurring in both standard WISE and WISE-S modes.

- **Shared Hallucination/Ambiguity:** In the "Concept Car" example, both $\tau_c$ and $\tau_d$ exhibit circular logic (tautology), failing to identify specific visual attributes. $\tau_c$ merely summarizes the vague reasoning of $\tau_d$.

**Conclusion:** These failure cases powerfully demonstrate the effectiveness of our **Self-Distillation** objective. The model successfully internalized the reasoning—whether strong or weak—into a compressed form. The errors stem from the backbone model's capabilities, not the thought compression process itself. This confirms that WISE-S provides a "lossless" speedup in terms of reasoning quality.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084



(a) **Qualitative Comparison: Efficiency vs. Distraction.** The figure demonstrates how WISE-S acts as a sufficient statistic, contrasting its focused reasoning with the baseline's verbose failure. **(Left) Functional Reasoning:** While WISE-S ($\tau_c$) efficiently identifies the glass by its affordance ("take small sips"), the baseline **Seg-Zero** (bottom) suffers from *reasoning drift*. It over-analyzes the prompt constraints, leading to a hallucinated conclusion about a non-existent "small textured object" to the right. This vivid example illustrates how thought compression can prevent the model from getting lost in irrelevant visual details. **(Right) Identity Grounding:** WISE-S correctly filters out spatial redundancy to focus on the discriminative visual attribute ("license plate") required to establish identity.



(b) Visual Attribute Grounding Case

Figure 1: **Success Case Study Examples.** The figures demonstrate how WISE-S acts as a sufficient statistic. (a) Shows the model efficiently identifying an object by its function ("take small sips") without verbose description. (b) Shows the model correctly filtering out spatial redundancy to focus on discriminative visual attributes ("red and white obstacle").

(a) Grounding Limitation Case



(b) Circular Logic Case

Figure 2: **Analysis of Failure Cases: Semantic Fidelity amidst Grounding Errors.** These examples illustrate instances where the model fails to produce an accurate mask (Low IoU). Crucially, however, the Concise Rationale ($\tau_c$) remains a **faithful summary** of the Detailed Explanation ($\tau_d$). (a) In the top example, both $\tau_c$ and $\tau_d$ correctly identify the semantic topic but struggle to ground the abstract concept to specific pixels. (b) In the bottom example, both rationales exhibit circular logic without identifying distinct visual attributes. This indicates that the failures stem from the **underlying limitations** of the base model's spatial grounding capabilities or reasoning loops, rather than information loss caused by the WISE-S compression mechanism.