# Exploring the Limits of Semantic Image Compression at Micro-Bits Per Pixel

Jordan Dotzel[1], Bahaa Kotb[11], James Dotzel[2], Mohamed Abdelfattah[1], Zhiru Zhang[1]

[1]Cornell University, [2]Penn State University
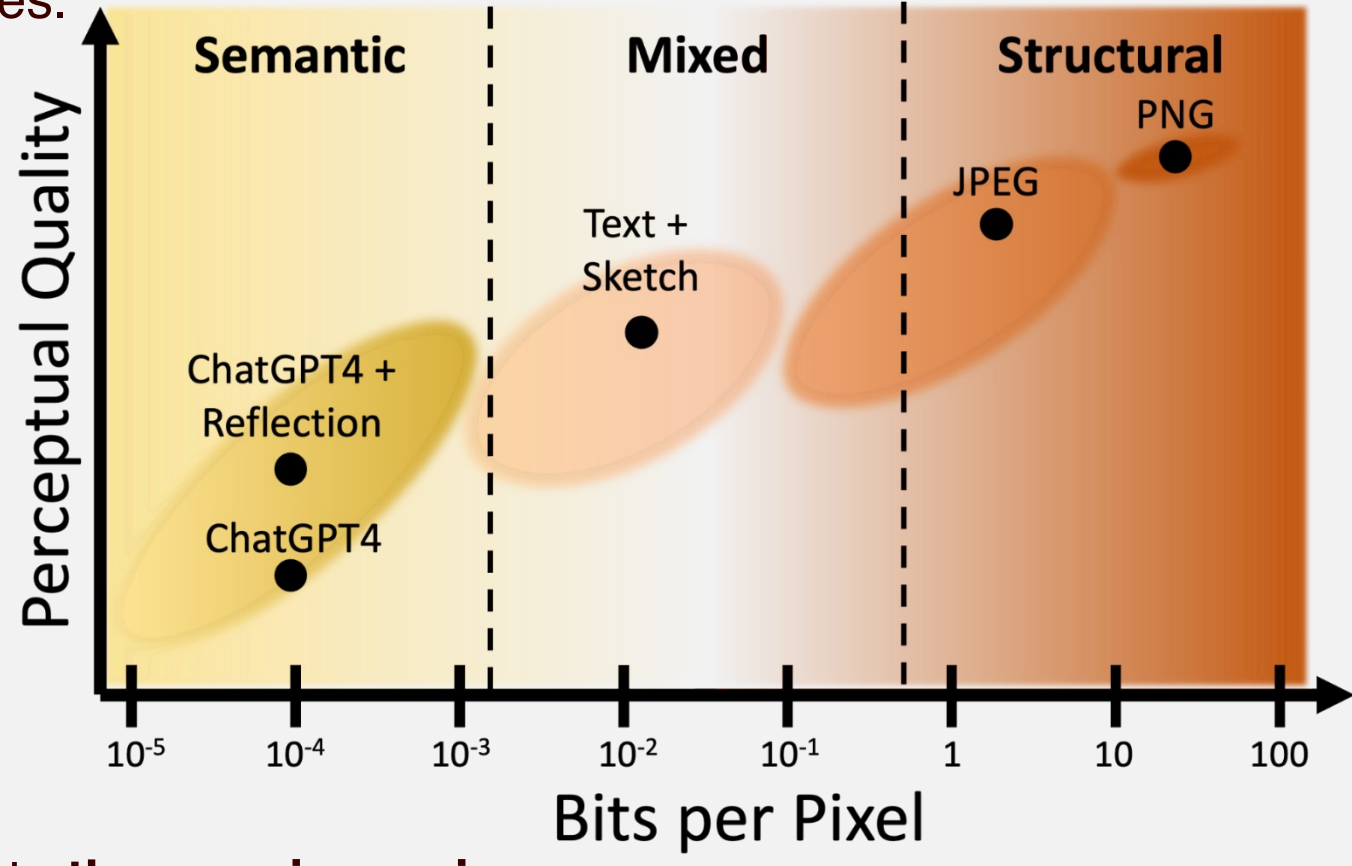
Equal Contribution

Tiny Paper at ICLR 2024

## Background

**Modern image compression** types:
- **Lossless**: Retains all information.
- **Lossy**: Loses some information, transferring data at **lower bitrates** to remove imperceptible details. For instance, **JPEG** retains only key human-perceptible frequencies.



Classification into **three major regions**:
- **Structural region**: Keeps original pixel structure, with JPEG operating between $10^{-1}$ to 10 bpp
- **Semantic region:** Focuses on essential human-centric information, achieving compression down to micro-bits per pixel ($\mu$bpp).
- **Mixed region**: Merges semantic and structural details, like **Text + Sketch**, which operates at $10^{-3}$ to $10^{-2}$ bpp.
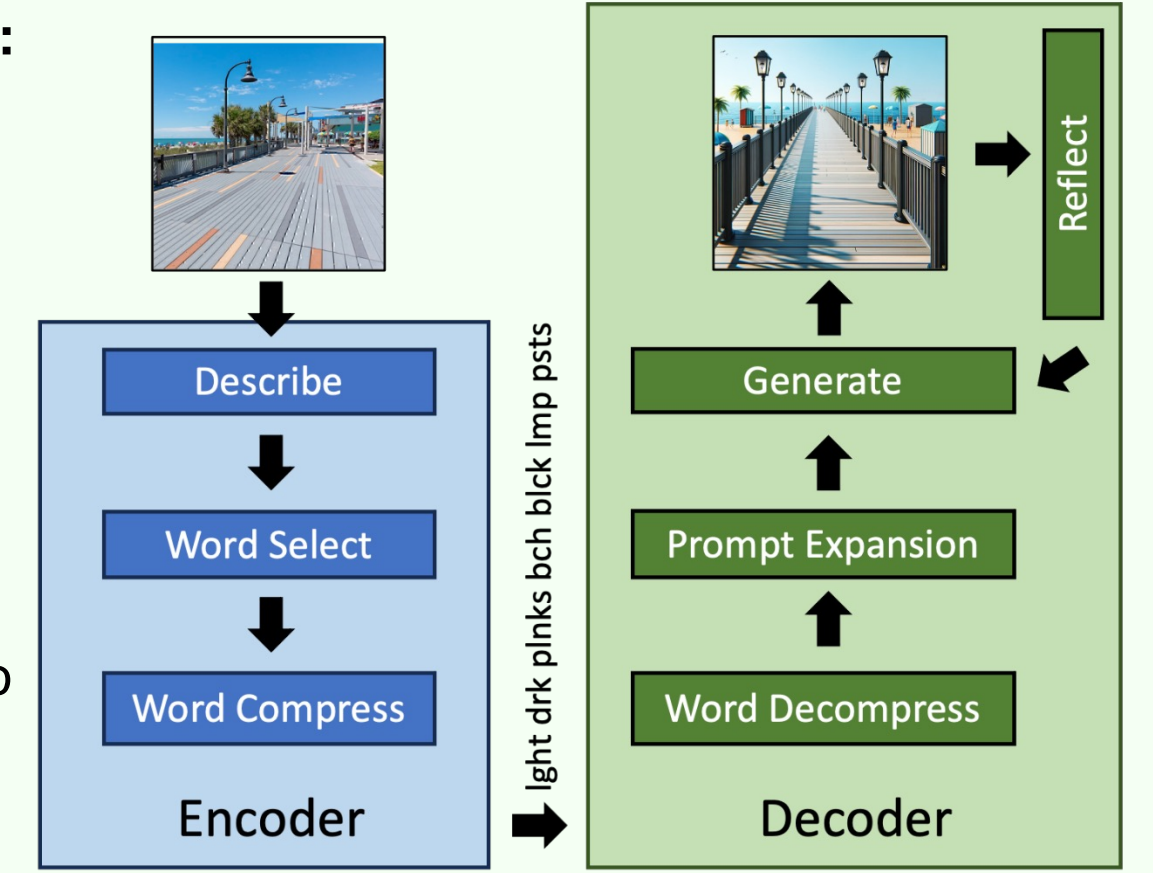
## Compression Methodology

**Base Models:** Uses **GPT-4V** as encoder and **DALL-E3** as decoder for advanced language compression and image generation.

**Encoding Process:**

- Analyzes and extracts key details from images, selecting the most important words and compressing them into a **four-bit representation** using only fifteen common consonants.

- **GPT-4V** optimizes text by substituting synonyms and minimizing characters.

**Decoding and Image Generation:**

- Compressed text is expanded into a natural prompt for **DALL-E3** to generate an image.

- Utilizes an **image reflection technique** for iterative enhancement if the context exceeds **500 $\mu$bpp**, typically requiring one or two iterations to refine the image.



## Compression Examples



Compression examples demonstrate a **progressive loss of contextual details** like room color and figure position, yet landmarks like the Taj Mahal retain **significant information at extremely low bitrates**, revealing a focus on arbitrary nuanced data like jacket color heavy compression.

## Reflection Examples



## Experimental Results



ChatGPT4 is a proof-of-concept for useful semantic compression, outperforming DALL-E3's generative abilities. The **practical limit of 100 $\mu$bpp was identified**, demonstrating that semantic compression can achieve significant size reduction with manageable loss in quality. The **reflection technique was crucial** for improving fidelity in the reconstructed images.