# Cheat Sheet: Foundations of Multimodal AI

| Concept | Description | Implementation Example |
|---|---|---|
| **Image captioning** | Generate descriptive captions for images using multimodal models. | ```python\ndef generate_image_caption(model, encoded_image):\n    """Generate a descriptive caption for an image."""\n    prompt = "Please provide a detailed description of this image."\n    return send_multimodal_query(model, encoded_image, prompt)\n# Example usage\ncaption = generate_image_caption(model, encoded_image)\nprint("Image Caption:", caption)\n``` |
| **Image processing** | Basic image processing and encoding for multimodal applications. | ```python\nimport base64\nfrom PIL import Image\nfrom io import BytesIO\ndef encode_image(image_path):\n    """Convert image to base64 for model input."""\n    with open(image_path, "rb") as image_file:\n        encoded_string = base64.b64encode(image_file.read()).decode('utf-8')\n    return encoded_string\ndef process_image(image_path, target_size=(224, 224)):\n    """Process image for model input."""\n    image = Image.open(image_path)\n    image = image.resize(target_size)\n    return image\n# Example usage\nimage_path = "example.jpg"\nencoded_image = encode_image(image_path)\nprocessed_image = process_image(image_path)\n``` |
| **Multimodal model setup** | Basic setup for working with multimodal AI models using IBM watsonx.ai platform. | ```python\nfrom ibm_watsonx_ai import Credentials\nfrom ibm_watsonx_ai.foundation_models import ModelInference\nfrom ibm_watsonx_ai.foundation_models.schema import TextChatParameters\ncredentials = Credentials(\n    url="https://us-south.ml.cloud.ibm.com",\n)\nparams = TextChatParameters(\n    temperature=0.2,\n    top_p=0.5,\n    max_tokens=2000\n)\nmodel = ModelInference(\n    model_id="meta-llama/llama-3-2-90b-vision-instruct",\n    credentials=credentials,\n    project_id="skills-network",\n    params=params\n)\n``` |
| **Multimodal query** | Send a combined text and image query to a multimodal model. | ```python\ndef send_multimodal_query(model, encoded_image, prompt):\n    """Send combined text and image query to model."""\n    messages = [{\n        "role": "user",\n        "content": [\n            {\n                "type": "text",\n                "text": prompt\n            },\n            {\n                "type": "image_url",\n                "image_url": {\n                    "url": f"data:image/jpeg;base64,{encoded_image}"\n                }\n            }\n        ]\n    }]\n    response = model.chat(messages=messages)\n    return response['choices'][0]['message']['content']\n# Example usage\nresponse = send_multimodal_query(model, encoded_image, prompt)\nprint("Model Response:", response)\n``` |

| | | |
|---|---|---|
| **Text processing** | Basic text processing and prompt engineering for multimodal applications. | ```python
def create_prompt(image_description, user_query):
    """Create a structured prompt for multimodal analysis."""
    prompt = f"""
Analyze the following image and answer the question.
Image Description: {image_description}
User Query: {user_query}
Please provide a detailed response that:
1. Describes the image content
2. Answers the specific question
3. Provides relevant context
"""
    return prompt
# Example usage
image_desc = "A cat sitting on a windowsill"
user_question = "What is the cat doing?"
prompt = create_prompt(image_desc, user_question)
``` |
| **Visual question answering** | Answer questions about image content using multimodal models. | ```python
def visual_question_answering(model, encoded_image, question):
    """Answer questions about image content."""
    prompt = f"Please answer the following question about the image: {question}"
    return send_multimodal_query(model, encoded_image, prompt)
# Example usage
question = "What color is the cat in the image?"
answer = visual_question_answering(model, encoded_image, question)
print("Answer:", answer)
``` |

# Author

[Ricky Shi](#)