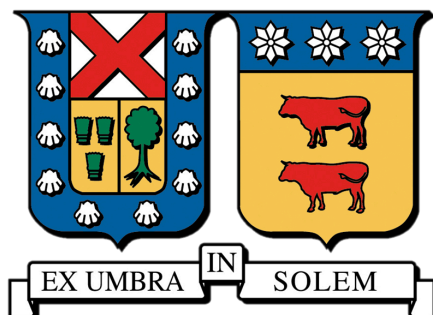


UNIVERSIDAD TÉCNICA FEDERICO SANTA MARÍA
DEPARTAMENTO DE ELECTRÓNICA
VALPARAÍSO – CHILE



Diseño y Desarrollo de Plataforma de Apreciación Artística para Personas con Discapacidad Visual

Alexey Nikolay Mitjaew Hupat

Memoria de Titulación para optar al título de Ingeniero Civil Telemático

Profesor Guía:

Patricio Olivares Roncagliolo

Profesor Correferente:

Nicolás Torres Rudloff

Diciembre 2025

Agradecimientos

A mi familia,
mi perro Elvis
y el buen café.

Resumen

Este documento expone el desarrollo de una plataforma accesible de apreciación artística para público con discapacidad visual. La propuesta consiste en transformar obras en representaciones sonoras multifacéticas, que integran descripciones narrativas detalladas, contexto histórico-artístico y atmósferas auditivas generadas mediante modelos multimodales de lenguaje, sistemas de síntesis de voz de alta y modelos de generación sonora basada en texto.

La solución expuesta propone facilitar el acceso al arte mediante una experiencia adaptada, con el objetivo de ofrecer alternativas de apreciación sensorial para personas con discapacidad visual.

Palabras clave: *Accesibilidad / Arte inclusivo / Inteligencia artificial / Sonificación*

Lista de Figuras

| | |
|--|----|
| Figura 1 Multimodal LLMs | 10 |
| Figura 2 Tacotron 2 Architecture | 11 |
| Figura 3 AudioLDM | 12 |
| Figura 4 Arquitectura General de Sistema | 13 |
| Figura 5 Flujo de Navegación Frontend | 15 |

Lista de Tablas

Contenidos

| | |
|--|----|
| Introducción | 6 |
| Objetivo General | 7 |
| Objetivos Particulares | 7 |
| Estructura | 7 |
| Marco Teórico | 8 |
| Modelos de Lenguaje Multimodales | 8 |
| Modelos Generativos de Audio | 8 |
| Estándares de Accesibilidad | 9 |
| Estado del Arte | 10 |
| Modelos de Lenguaje Multimodales | 10 |
| Modelos de Texto a Voz | 11 |
| Modelos Generativos de Ambiente Sonoro | 12 |
| Desarrollo de la Plataforma | 13 |
| API Gateway | 14 |
| Narrador de Contexto | 14 |
| Generador de Audio Descriptivo | 14 |
| Generador de Sonidos Ambientales | 14 |
| Catálogo de Obras/Imagenes | 15 |
| Frontend | 15 |
| Limitaciones | 15 |
| Resultados | 16 |
| Conclusiones | 16 |
| Bibliografía | 16 |
| Anexo | 17 |
| Prompts | 17 |
| Utilidades Generales | 17 |
| Narrador de Contexto | 24 |

| | |
|--|----|
| Generación de Audio Descriptivo | 26 |
| Generación de Sonidos Ambientales | 28 |
| Procesamiento de Catálogo de Obras | 30 |
| Text to Speech | 30 |
| Modelos Generativos de Audio Ambiental | 31 |
| AudioLDM | 31 |

Introducción

El acceso equitativo al arte constituye un componente esencial del desarrollo cultural y social, sin embargo, las personas con discapacidad visual continúan enfrentando barreras significativas para disfrutar de obras visuales en condiciones comparables al resto del público. Aunque las tecnologías de apoyo han avanzado, aún existe una brecha entre la experiencia estética que ofrecen los museos y plataformas digitales tradicionales y las necesidades sensoriales de quienes no pueden percibir elementos visuales directamente.

En este contexto, la inteligencia artificial emerge como una oportunidad para reimaginar la forma en que se transmite el contenido artístico, permitiendo crear experiencias inmersivas que integren narrativa, sonido y contextualización cultural. Este proyecto desarrolla una plataforma que transforma obras de arte en representaciones sonoras enriquecidas, combinando modelos multimodales capaces de interpretar imágenes, sistemas de síntesis de voz de alta naturalidad y tecnologías generativas orientadas a la creación de paisajes sonoros coherentes con la obra original.

La presente investigación se enmarca en los principios de accesibilidad universal y diseño inclusivo, procurando que la experiencia resultante no solo sea funcional, sino también significativa desde una perspectiva estética. La plataforma busca responder a la necesidad de ofrecer alternativas sensoriales que amplíen el acceso al patrimonio artístico y cultural, fortaleciendo la inclusión mediante herramientas tecnológicas avanzadas.

Objetivo General

Mejorar la experiencia estética y promover la accesibilidad universal a las obras de arte para personas con discapacidad visual mediante el diseño e implementación de una plataforma accesible que utilice inteligencia artificial para transformar automáticamente imágenes de obras artísticas en paisajes sonoros narrativos, contextuales y ambientales.

Objetivos Particulares

- **Proveer un catálogo curado de obras de dominio público**, basado en fuentes autoritativas y verificadas, que permita a los usuarios acceder a contenido confiable, culturalmente riguroso y legalmente seguro.
- **Generar descripciones auditivas de alta calidad** utilizando modelos de lenguaje multimodales capaces de interpretar intuitivamente los elementos visuales presentes en cada obra, garantizando una representación objetiva y comprensible de su contenido.
- **Crear narraciones históricas y biográficas** fundamentadas en fuentes reconocidas, con el fin de contextualizar la obra, su autor(a) y su relevancia artística dentro de un marco cultural accesible.
- **Producir ambientes sonoros inmersivos** que recreen atmósferas coherentes con la escena o época sugerida por la obra, mediante modelos generativos capaces de transformar imágenes en audio evocativo y sensorialmente enriquecido.
- **Diseñar una interfaz centrada en accesibilidad**, que permita la navegación autónoma de personas en situación de discapacidad visual, incorporando criterios y normativas vigentes en materia de accesibilidad universal y diseño inclusivo.

Estructura

El presente análisis incluye:

- *Revisión del estado del arte sobre modelos generativos y estándares de accesibilidad.*
- *Evaluación del marco teórico utilizado para el desarrollo de la plataforma.*
- *Descripción de la plataforma desarrollada a nivel de componentes.*
- *Revisión de resultados obtenidos y conclusiones al respecto.*

Marco Teórico

Modelos de Lenguaje Multimodales

Los modelos de lenguaje (LLMs, por sus siglas en inglés) representan sistemas computacionales diseñados para procesar, comprender y generar texto en lenguaje natural mediante arquitecturas basadas en aprendizaje profundo. Estos modelos, entrenados en vastos corpus textuales, capturan patrones lingüísticos, semánticos y contextuales, permitiendo tareas como la traducción automática, la generación de respuestas coherentes o el análisis de sentimientos. Su evolución ha dado paso a los modelos de lenguaje multimodales, los cuales extienden estas capacidades al integrar múltiples formas de información, como imágenes y audio.

Estos modelos multimodales no solo interpretan contenido visual o auditivo, sino que también generan descripciones detalladas de escenas, objetos o composiciones artísticas, traduciendo elementos visuales en narrativas estructuradas. Además, su capacidad multimodal les permite procesar entradas de voz y participar en flujos conversacionales voz a voz, facilitando interacciones más naturales y accesibles. Esta convergencia de modalidades (texto, imagen y sonido) facilita la creación de sistemas más ricos y adaptables, capaces de ofrecer interpretaciones contextualizadas y personalizadas según las necesidades de percepción del usuario.

Modelos Generativos de Audio

Los modelos generativos de audio representan sistemas capaces de crear sonidos, música y voces artificiales a partir de descripciones semánticas o patrones aprendidos. Estos modelos, basados en arquitecturas de aprendizaje profundo, pueden sintetizar desde elementos literales —como el canto de pájaros, el sonido de una tormenta o pasos sobre grava— hasta composiciones musicales complejas, ajustándose a estilos, géneros o emociones específicas. Su funcionamiento se sustenta en la interpretación de instrucciones textuales o parámetros acústicos, permitiendo generar audio coherente y contextualizado

sin necesidad de muestras preexistentes. Esta capacidad no solo amplía las posibilidades creativas en producción musical y diseño sonoro, sino que también facilita la creación de entornos auditivos personalizados, como paisajes sonoros para aplicaciones de accesibilidad o asistentes de voz con tonos y matices más naturales.

Estándares de Accesibilidad

Los estándares de accesibilidad proporcionan el marco normativo que guía el diseño de interfaces inclusivas, asegurando que los contenidos digitales sean utilizables por personas con diversas capacidades sensoriales, motoras o cognitivas. Entre los más relevantes se encuentran las Pautas de Accesibilidad para el Contenido Web (WCAG), que establecen criterios relacionados con la perceptibilidad, operabilidad, comprensibilidad y robustez de los sistemas web. En el contexto chileno, estas directrices se integran mediante la Ley N.º 20.422, que regula la igualdad de oportunidades e incorpora exigencias específicas para tecnologías asistivas. La plataforma desarrollada en este proyecto adopta estos lineamientos para garantizar compatibilidad con lectores de pantalla, navegación por teclado, descripciones alternativas, controles auditivos accesibles y una estructura de interacción que facilite la exploración autónoma.

El cumplimiento de estos estándares es esencial para asegurar que la solución propuesta no solo sea técnicamente avanzada, sino también verdaderamente inclusiva y centrada en las necesidades de la comunidad usuaria.

Estado del Arte

Modelos de Lenguaje Multimodales

Durante los últimos años los modelos multimodales de lenguaje han experimentado avances significativos, integrando capacidades de procesamiento de texto, imagen, audio y video. Este progreso ha sido impulsado por arquitecturas unificadas basadas en Transformers que permiten compartir un espacio de representación común entre modalidades. Avances como los *Vision-Language Encoders (VLMs)*, los *Audio-Text aligners* y la *Tokenización Unificada* han permitido integrar imágenes, audio y texto como secuencias compatibles. Además, técnicas como el *cross-attention multimodal*, la *late fusion optimizada* y el *modality routing* han mejorado la coherencia entre modalidades, habilitando tareas complejas de razonamiento y generación conjunta en múltiples formatos.

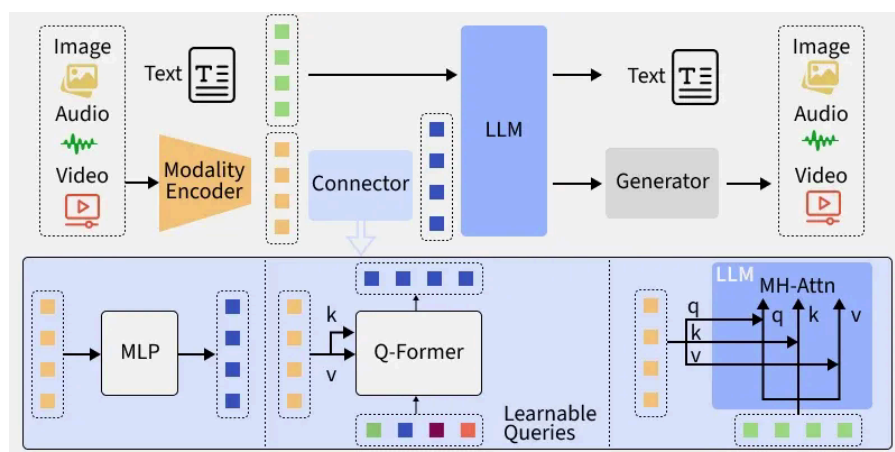


Figura 1: Multimodal LLMs

Entre los avances destacados hasta la fecha (2025) se encuentran:

- **Llama v4 Maverick (Meta):** Modelo multimodal con descripción semántica detallada de imágenes y análisis de escenas. Pesos liberados para investigación.
- **GPT-4o (OpenAI) y Gemini 1.5 (Google):** Modelos con integración nativa de texto, imagen y audio. Sin pesos liberados.
- **Voice-to-Voice (V2V):** Sistemas como Whisper v3 (OpenAI) y SeamlessM4T (Meta) para transcripción y traducción en tiempo real. Voicebox (Meta) genera voz sintética con control de emociones. Whisper v3 tiene pesos liberados; SeamlessM4T y Voicebox no.

Modelos de Texto a Voz

En los últimos años, arquitecturas como Tacotron 2, VITS, Kokoro TTS y modelos basados en difusión/transformers han ganado relevancia. Estos modelos, generalmente ligeros en cantidad de parámetros, son razonables de autohostear con GPUs de grado consumidor e incluso en CPU en algunos casos. Generan voces naturales con control sobre prosodia, entonación y emociones. Kokoro TTS, de código abierto, destaca por su expresividad y ajuste de voces mediante manipulación del espacio latente, evitando reentrenamiento.

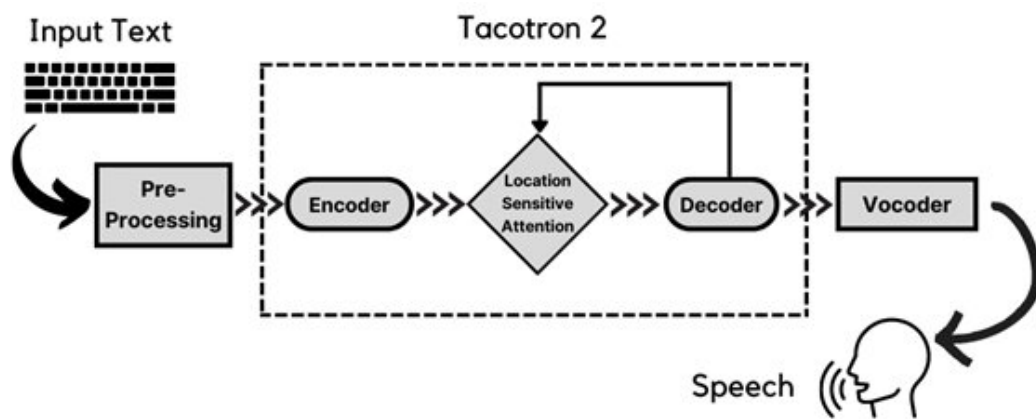


Figura 2: Tacotron 2 Architecture

Además de los modelos de código abierto, proveedores comerciales como Google Cloud, Amazon Polly, Microsoft Azure, ElevenLabs y OpenAI ofrecen soluciones TTS de alta calidad. Estas plataformas incorporan voces multilingües, opciones de personalización avanzada y capacidades de clonación vocal.

No obstante los avances observados, subsisten desafíos técnicos comunes en las alternativas investigadas, como la generación de expresividad controlable, la adaptación a entornos acústicos adversos y la síntesis eficiente de secuencias de larga duración.

Modelos Generativos de Ambiente Sonoro

La generación de ambientes sonoros de alta calidad sigue en fase experimental, sin soluciones comerciales consolidadas. Modelos como CLAP permiten mapear texto y audio en un espacio semántico compartido, aunque no generan audio directamente.

Entre los modelos evaluados para la plataforma se encuentran:

- **I Hear Your True Colors:** La arquitectura propuesta en este paper combina el uso de un VQVAE, transformers y CLIP. El VQVAE extrae representaciones jerárquicas como secuencias discretas, los transformers las modelan de forma autorregresiva y CLIP alinea el audio con lo visual. Aunque el código es abierto, no existen fuentes públicas de pesos preentrenados. Sin versiones comerciales.
- **AudioLDM:** Utiliza difusión para generar espectrogramas condicionados por texto, incorporando embeddings de CLAP para mejorar la alineación semántica entre descripciones textuales y contenido sonoro. Presenta resultados prometedores en coherencia semántica y diversidad acústica. El modelo preentrenado se puede encontrar con pesos abiertos.

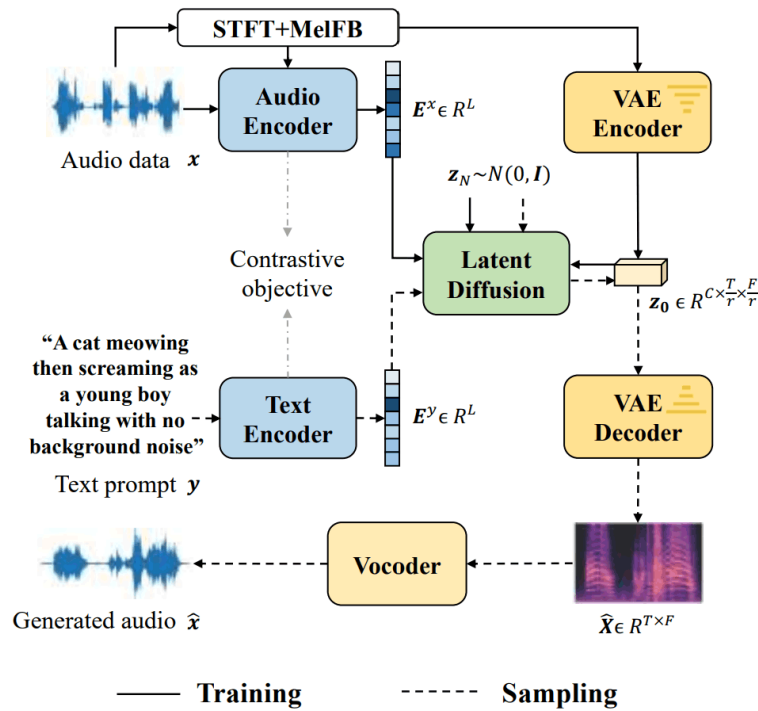


Figura 3: AudioLDM

Desarrollo de la Plataforma

La arquitectura propuesta sigue un modelo cliente-servidor e incorpora un *API Gateway* como componente central. Este elemento cumple dos funciones principales: por un lado, centraliza las operaciones computacionalmente intensivas (principalmente inferencia de modelos de aprendizaje automático y las solicitudes a APIs externas), y por otro, gestiona un catálogo centralizado de recursos. Además, se desarrolla un *frontend* que aplica principios de accesibilidad web, con especial atención a las necesidades de usuarios con discapacidad visual.

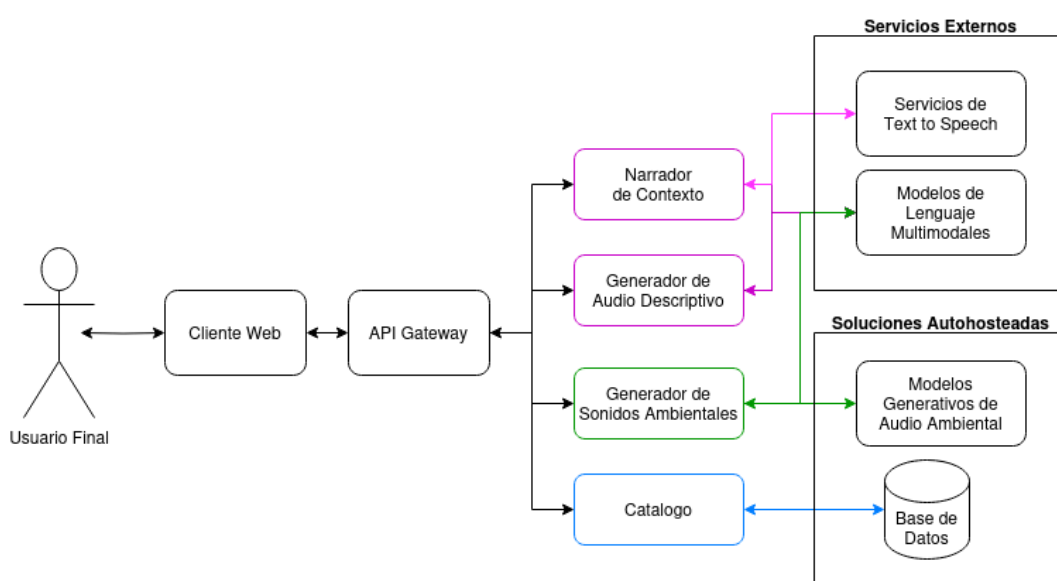


Figura 4: Arquitectura General de Sistema

API Gateway

La API Gateway consiste en un servidor de Django, encargado de controlar 4 servicios:

Narrador de Contexto

Este módulo se encarga de producir descripciones textuales detalladas y estructuradas de los contenidos visuales presentes en las obras. Para ello, emplea modelos de lenguaje multimodales, los cuales procesan tanto elementos gráficos como metadatos asociados mediante prompts especializados. Dichos prompts están diseñados para extraer información semántica relevante, garantizando una representación fiel y contextualizada del material original (**ANEXO PROMPT AQUI**). Posteriormente, las descripciones generadas son convertidas a formato de audio mediante sistemas de síntesis de voz (TTS) (**ANEXO TTS AQUI**), optimizados para ofrecer una experiencia auditiva clara y accesible.

Generador de Audio Descriptivo

Este módulo genera narraciones auditivas que contextualizan históricamente cada obra. Para ello, se recopila información de fuentes autoritativas (en este caso, se procesó manualmente un conjunto de datos con artículos de Wikipedia vinculados a cada pieza) (**ANEXO DATASET AQUI**). A partir de estos datos, se elabora un relato contextualizado y adaptado al caso de uso (**ANEXO PROMPT AQUI**), que posteriormente se convierte en audio mediante modelos de síntesis de voz (TTS) (**ANEXO TTS AQUI**), asegurando una experiencia auditiva coherente y accesible.

Generador de Sonidos Ambientales

Este módulo se encarga de la generación de ambientes sonoros que representan los contenidos literales (no abstractos) de sus imágenes de entrada. La heurística utilizada consiste en:

- Dividir el espacio en cuadrantes (en nuestro caso 9) (**ANEXO**)
- Procesar cada cuadrante con LLMs multimodales, obteniendo salidas estructuradas en JSON (**ANEXO**) representando cada elemento detectado.
- Por cada cuadrante generar una mezcla de sonido integrando todos los elementos detectados en cada respectiva sección. (**ANEXO**)

Catálogo de Obras/Imagenes

Se elaboró un conjunto de datos en formato JSON que incluye 30 obras artísticas, conteniendo todo el material requerido para el funcionamiento de los sistemas previamente descritos. Este dataset incorpora, además, enlaces a los artículos correspondientes de Wikipedia asociados a cada obra artística. **(ANEXO)**

Frontend

El frontend se implementó utilizando NextJS para proporcionar acceso al contenido gestionado por la API Gateway. En su diseño, se incorporaron textos alternativos y etiquetas HTML semánticas con el objetivo de garantizar una integración completa con herramientas de asistencia para usuarios con discapacidad visual. Asimismo, se desarrolló una interfaz de usuario que facilita la navegación mediante atajos de teclado, optimizando la accesibilidad y usabilidad del sistema.

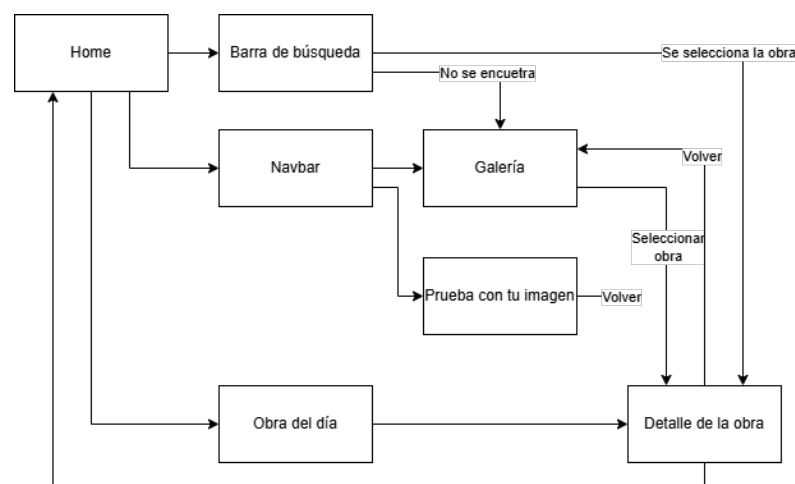


Figura 5: Flujo de Navegación Frontend

Limitaciones

Por limitaciones computacionales durante el desarrollo, los módulos no se integraron directamente en la API Gateway. En su lugar, se implementó un prototipo funcional en un cuadernillo Jupyter (<https://www.kaggle.com/code/alexeymitjaew/sound-generation-pipeline>) para validar el diseño y flujo de procesamiento. Esto permitió generar y preprocesar los conjuntos de datos necesarios, que fueron incorporados a la base de datos para su uso en la versión final.

Resultados

Conclusiones

[1] said [2] said

Bibliografía

- [1] Ö. Aksin *et al.*, «Effect of immobilization on catalytic characteristics of saturated Pd-N-heterocyclic carbenes in Mizoroki-Heck reactions», *J.~Organomet. Chem.*, vol. 691, n.º 13, pp. 3027-3036, 2006.
- [2] G. Westfahl, «The True Frontier». pp. 55-65.

Anexo

```
1 def test(x: str):  
2     print(f"String is: {x}")
```

Código 1: Some code

Como se ve en Código 1, podemos blablabla

Prompts

Utilidades Generales

```
1 TRANSLATOR_PROMPT = \ Python  
2 """Traduce el siguiente texto al español de manera directa y clara.  
3 Considera lo siguiente:  
4 - Mantén el significado exacto del texto original.  
5 - NO agregues explicaciones, comentarios o interpretaciones.  
6 - NO seas redundante.  
7 - Conserva el estilo y tono original del texto.  
8 - Escribe solo la traducción, sin títulos ni subtítulos.  
9  
10 Texto original: {TEXT}  
11 Traducción: """"  
12  
13 def translate(text):  
14     completion = client.chat.completions.create(  
15         #model="llama-3.1-8b-instant",  
16         model="openai/gpt-oss-120b",  
17         messages=[  
18             "role": "user",  
19             "content": [{  
20                 "type": "text",  
21                 "text": PARSE_PROMPT(TRANSLATOR_PROMPT, TEXT=text)  
22             }]  
23         ]],  
24         temperature=0,  
25         max_completion_tokens=1024,  
26         top_p=1,  
27         stream=False,  
28         stop=None,  
29     )  
30     content = completion.choices[0].message.content  
31     return content
```

Código 2: Prompt de Traducción

```

1 SOURCE_EXTRACTOR_PROMPT = \
2 """Eres un extractor de información precisa y concisa.
3 Tu tarea es leer el texto de un artículo y devolver una lista de
4 viñetas (•) con todos los hechos y datos relevantes,
5 escritos en español natural y claro.
6 Instrucciones:
7 1. Identifica los hechos principales, cifras, declaraciones, fechas,
8 nombres y conclusiones clave del texto.
9 2. No incluyas opiniones, lenguaje publicitario o frases
10 irrelevantes.
11 3. Si el texto está en inglés u otro idioma, traduce los puntos al
12 español correctamente.
13 4. Usa frases breves pero completas, cada una iniciando con un punto
14 (•).
15 5. Mantén el tono informativo y objetivo.
16
17 Texto del artículo: {TEXT}
18
19 Tu respuesta: """
20
21 def extract_text_from_url(url: str) -> str:
22     headers = {
23         "User-Agent": "Mozilla/5.0 (X11; Linux x86_64; rv:140.0)
24         Gecko/20100101 Firefox/140.0",
25     }
26     response = requests.get(url, headers=headers)
27     response.raise_for_status()
28
29     soup = BeautifulSoup(response.text, 'html.parser')
30
31     for tag in soup(['script', 'style']):
32         tag.decompose()
33
34     paragraphs = [p.get_text(strip=True) for p in soup.find_all('p')]
35     return '\n\n'.join(paragraphs)

```

Código 3: Extractor de Información

```
1  def extract_url_context(url):
2      text = extract_text_from_url(url)[:8_000]
3
4      completion = client.chat.completions.create(
5          model="openai/gpt-oss-120b",
6          messages=[{
7              "role": "user",
8              "content": [{
9                  "type": "text",
10                 "text": PARSE_PROMPT(SOURCE_EXTRACTOR_PROMPT,
11                                     TEXT=text)
12             }],
13             temperature=0,
14             max_completion_tokens=2048,
15             top_p=1,
16             stream=False,
17             stop=None,
18         )
19     content = completion.choices[0].message.content
20     return content
```

 Python

```

1  TECHNIQUE_PROMPT = \
2  """Explicame en detalle una técnica de pintura o dibujo.
3  Considera lo siguiente:
4  - Debe estar en un formato adecuado para narración.
5  - Debe ser conciso y corto en duración.
6  - NO digas en resumen.
7  - NO seas redundante.
8
9  - NO utilices títulos ni subtítulos, solo escribe en párrafos
10 narrados.
11
12 Técnica: {TECHNIQUE}
13 Explicación: """
14
15 def get_technique(technique):
16     completion = client.chat.completions.create(
17         #model="llama-3.1-8b-instant",
18         model="llama-3.3-70b-versatile",
19         messages=[{
20             "role": "user",
21             "content": [{
22                 "type": "text",
23                 "text": PARSE_PROMPT(TECHNIQUE_PROMPT,
24                                     TECHNIQUE=technique)
25             }]
26         }],
27         temperature=0,
28         max_completion_tokens=1024,
29         top_p=1,
30         stream=False,
31         stop=None,
32     )
33     content = completion.choices[0].message.content
34     return content

```

Código 5: Generador de Texto de Técnica

```

1  PERIOD_PROMPT = \
2  """Explicame en detalle un estilo pictórico en particular.
3  Considera lo siguiente:
4  - Debe estar en un formato adecuado para narración.
5  - Debe ser conciso y corto en duración.
6  - NO digas en resumen.
7  - NO seas redundante.
8
9  - NO utilices títulos ni subtítulos, solo escribe en párrafos
10 narrados.
11
12 Estilo pictórico: {PERIOD}
13 Explicación: """
14
15 def get_period(period):
16     completion = client.chat.completions.create(
17         #model="llama-3.1-8b-instant",
18         model="llama-3.3-70b-versatile",
19         messages=[{
20             "role": "user",
21             "content": [{
22                 "type": "text",
23                 "text": PARSE_PROMPT(PERIOD_PROMPT, PERIOD=period)
24             }]
25         }],
26         temperature=0,
27         max_completion_tokens=1024,
28         top_p=1,
29         stream=False,
30         stop=None,
31     )
32     content = completion.choices[0].message.content
33     return content

```

Código 6: Generador de Texto de Periodo

```

1  AUTHOR_BIO_PROMPT = \
2  """Eres un experto en redacción biográfica y narrativa.
   Tu tarea es leer el siguiente texto y redactar una biografía breve y
3  fluida del autor o artista mencionado, adecuada para acompañar una
   ficha de obra de arte.
4
5  Instrucciones:
   1. Resume los aspectos esenciales de la vida y obra del autor:
6   formación, estilo, periodo histórico, aportes y relevancia artística.
   2. Mantén un tono natural, informativo y narrativo, sin usar listas
7   ni formato enciclopédico.
   3. Evita redundancias, tecnicismos innecesarios y frases genéricas.
8   4. No incluyas fechas o datos no presentes en el texto fuente.
9   5. No menciones el proceso de redacción ni expresiones como "esta
10  biografía trata sobre".
   6. Si el texto está en otro idioma, tradúcelo al español de manera
11  natural.
12  7. Extensión máxima: dos párrafos.
13  8. Evita completamente el uso de caracteres especiales
14
15  Texto: {TEXT}
16
17  Tu respuesta: """
18
19  def get_bio(url):
20      text = extract_url_context(url)
21      completion = client.chat.completions.create(
22          model="llama-3.3-70b-versatile",
23          messages=[{
24              "role": "user",
25              "content": [{
26                  "type": "text",
27                  "text": PARSE_PROMPT(AUTHOR_BIO_PROMPT, TEXT=text)
28              }]
29          }],
30          temperature=0,
31          max_completion_tokens=1024,
32          top_p=1,
33          stream=False,
34          stop=None,
35      )
36      content = completion.choices[0].message.content
37      return content

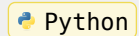
```

Código 7: Author Bio Generation

```

1  def resize_if_oversized(image: Image.Image, max_dim: int =
2      1024) -> Image.Image:
3      width, height = image.size
4      if max(width, height) <= max_dim:
5          return image # nothing to do
6
7      # scale preserving proportions
8      scale = max_dim / max(width, height)
9      new_size = (int(width * scale), int(height * scale))
10     return image.resize(new_size, Image.Resampling.LANCZOS)
11
12 def download_img(url, max_dim: int = 1024):
13     headers = { "User-Agent": "Chrome/51.0.2704.103 Safari/537.36" }
14     r = requests.get(url, headers=headers)
15     im = Image.open(BytesIO(r.content))
16     im = resize_if_oversized(im, max_dim)
17     return im
18
19 def image_to_base64(image: Image.Image, format="PNG") -> str:
20     """Convert PIL Image to base64 string."""
21     buffered = BytesIO()
22     image.save(buffered, format=format)
23     return base64.b64encode(buffered.getvalue()).decode("utf-8")

```



Código 8: Image Parsing

Narrador de Contexto

Python

```

1 SOURCE_SUMMARY_PROMPT = \
2 """Eres un experto en redacción narrativa.
   Tu tarea es leer el siguiente texto y escribir un resumen breve
3 (máximo dos párrafos) en un estilo natural, fluido y adecuado para
   narración oral o escrita.
4
5 Instrucciones:
6     1. Mantén un tono narrativo, como si estuvieras contando los hechos
7       de manera clara y atractiva.
8     2. Sé conciso: evita repeticiones, rodeos o frases innecesarias.
9     3. No incluyas títulos, subtítulos ni listas.
10    4. No uses expresiones como “en resumen”, “este artículo trata
11       sobre”, ni menciones al proceso de resumen.
12    5. Si el texto está en otro idioma, tradúcelo al español con
13       naturalidad.
14    6. Enfócate en los hechos principales y el contexto historico de la
15       producción de la obra.
16
17 Texto: {TEXT}
18
19 Tu respuesta: """
20
21 def get_narration(url):
22     text = extract_url_context(url)
23     completion = client.chat.completions.create(
24         model="llama-3.3-70b-versatile",
25         messages=[{
26             "role": "user",
27             "content": [{
28                 "type": "text",
29                 "text": PARSE_PROMPT(SOURCE_SUMMARY_PROMPT,
30                                     TEXT=text)
31             }]
32         }],
33         temperature=0,
34         max_completion_tokens=2048,
35         top_p=1,
36         stream=False,
37         stop=None,
38     )
39     content = completion.choices[0].message.content
40     return content

```

Generación de Audio Descriptivo

Python

```

1 DESCRIPTION_PROMPT = \
2     """Describe los elementos retratados en la imagen.
3     Considera lo siguiente:
4     - Debe estar en un formato adecuado para narración.
5     - Debe ser conciso y corto en duración.
6     - NO digas en resumen.
7     - NO seas redundante.
8     - NO utilices títulos ni subtítulos, solo escribe en párrafos
9     narrados.
10
11 Elementos retratados: """
12
13 def get_description(url):
14     image = download_img(url)
15     b64_image = image_to_base64(image)
16     completion = client.chat.completions.create(
17         model="meta-llama/llama-4-scout-17b-16e-instruct",
18         messages=[
19             {
20                 "role": "user",
21                 "content": [
22                     {
23                         "type": "text",
24                         "text": DESCRIPTION_PROMPT
25                     },
26                     {
27                         "type": "image_url",
28                         "image_url": {
29                             "url": f"data:image/png;base64,
30                                 {b64_image}"
31                         }
32                     }
33                 ],
34                 temperature=0,
35                 max_completion_tokens=1024,
36                 top_p=1,
37                 stream=False,
38                 stop=None,
39             )
40     return completion.choices[0].message.content

```

Código 10: Generación de texto para audios descriptivos

Generación de Sonidos Ambientales

```
1  def get_quadrants(image: Image, N:int=3):  
2      total_width, total_height = image.size  
3      w, h = total_width//N, total_height//N  
4      w_pos = [(w*i, w*(i+1)) for i in range(N)]  
5      h_pos = [(h*i, h*(i+1)) for i in range(N)]  
6      boxes = [  
7          (w_pos[i][0], h_pos[j][0], w_pos[i][1], h_pos[j][1])  
8          for j in range(N) for i in range(N)  
9      ]  
10     cropped = [  
11         image.crop(box)  
12         for box in boxes  
13     ]  
14     coords= [  
15         (  
16             (i/N + (i+1)/N)/2,  
17             (j/N + (j+1)/N)/2,  
18         )  
19         for j in range(N) for i in range(N)  
20     ]  
21  
22     return cropped, coords
```

Código 11: Cropping de imagen en cuadrantes

```

1  EXTRACTOR_PROMPT = (
2      """Extract audio ambience generation prompts from the following
3      image.
4      Consider the following:
5      - Your output should be a list of json objects containing the keys:
6        {"is_background": bool, "object": string}
7      - Only indicate objects with sound ambience relevance, ignore muted
8        elements
9      - Ignore abstract elements
10     - Indicate 3 elements at most
11     - Group elements when there are more than 1 depicted
12     - You should only output ONE list
13     - Alive elements should be added with their corresponding sound
14       description (bear growling, dog barking)
15     - Any other non alive element should just be specified as the element
16       itself
17     - Output just the required JSON results
18     JSON RESULT: """)
19
20 def get_semantic_elements(image: Image, max_attempts=3):
21     b64_image = image_to_base64(image)
22     for _ in range(max_attempts):
23         try:
24             completion = client.chat.completions.create(
25                 model="meta-llama/llama-4-scout-17b-16e-instruct",
26                 messages=[
27                     {
28                         "role": "user",
29                         "content": [
30                             {
31                                 "type": "text",
32                                 "text": EXTRACTOR_PROMPT
33                             },
34                             {
35                                 "type": "image_url",
36                                 "image_url": {
37                                     "url": f"data:image/png;base64,
38                                         {b64_image}"
39                                 }
40                             }
41                         ]
42                     }
43                 ]
44             )
45             return [{"object": obj["object"], "is_background": obj["is_background"]}
46                     for obj in completion.choices[0].message.content]

```

Código 12: Detección de elementos en la obra

Procesamiento de Catálogo de Obras

Text to Speech

```
1 # !pip install -q kokoro>=0.9.4 soundfile  Python
2 # !apt-get -qq -y install espeak-ng > /dev/null 2>&1
3
4 from kokoro import KPipeline
5 from IPython.display import display, Audio
6 import soundfile as sf
7 import torch
8
9 tts_pipeline = KPipeline(lang_code='e')python
```

Código 13: Setup Kokoro TTS

```
1 def tts(text, filename):  Python
2     generator = tts_pipeline(
3         text,
4         voice='em_santa',
5         speed=1, split_pattern=r'\n+'
6     )
7
8     pause_duration = 0.5
9     silence = np.zeros(int(24000 * pause_duration), dtype=np.float32)
10    all_audio = []
11
12    for i, (gs, ps, audio) in enumerate(generator):
13        all_audio.append(audio)
14        all_audio.append(silence)
15
16    final_audio = np.concatenate(all_audio)
17    sf.write(f"{filename}.wav", final_audio, 24000)
```

Código 14: Generación de audio con Kokoro TTS

Modelos Generativos de Audio Ambiental

AudioLDM

```
1 from IPython.display import Audio
2 from diffusers import AudioLDMPipeline
3 import torch
4
5 repo_id = "cvssp/audioldm"
6 pipe = AudioLDMPipeline.from_pretrained(repo_id,
7 torch_dtype=torch.float16)
8 pipe = pipe.to("cuda")
9 generator = torch.Generator("cuda").manual_seed(0)
```

Código 15: Audio LDM Setup

```

1  from pydub import AudioSegment
2  import torch
3  import numpy as np
4  import tempfile
5  import os
6
7  def tensor_to_audiosegment(audio_tensor, sample_rate=16000):
8      samples = (audio_tensor * 32767).astype(np.int16)
9      return AudioSegment(
10         samples.tobytes(),
11         frame_rate=sample_rate,
12         sample_width=2, # 16-bit = 2 bytes
13         channels=1
14     )
15
16  VOLUME_DB_BACKGROUND = -20
17  VOLUME_DB_FOREGROUND = -5
18  generator = torch.Generator("cuda").manual_seed(0)
19
20  def generate_ambient_sound(scene, filename):
21      final_mix = AudioSegment.silent(duration=20000)
22
23      for i, obj in enumerate(scene):
24          prompt = obj['object']
25          negative_prompt = ["low quality, average quality"]
26
27          result = pipe(
28              [prompt],
29              negative_prompt=negative_prompt,
30              num_inference_steps=30,
31              audio_length_in_s=20,
32              generator=generator,
33              return_dict=True
34          )
35
36          audio_tensor = result.audios[0].squeeze()
37          audio_segment = tensor_to_audiosegment(audio_tensor)
38          volume_db = VOLUME_DB_BACKGROUND if obj["is_background"] else
39              VOLUME_DB_FOREGROUND
40          audio_segment = audio_segment + volume_db
41          final_mix = final_mix.overlay(audio_segment)
42
43      final_mix.export(f"{filename}.wav", format="wav")

```

Código 16: Generación de mezcla de sonido ambiental con AudioLDM