

INGENIERÍA EN SONIDO

Título de Tesis

Subtítulo de la tesis (si lo tuviera)

*Tesis final presentada para obtener el título de Ingeniero de Sonido de la
Universidad Nacional de Tres de Febrero (UNTREF)*

TESISTA: Matías Di Bernardo (42.229.438)
TUTOR/A: Nombre y apellido (Ing., PhD., etc.)
COTUTOR/A: Gala Lucía Gonzalez Barrios (Lic.)

AGRADECIMIENTOS

Se propone incluir este apartado, donde se debe agradecer primeramente a las autoridades de la Universidad, al coordinador de la carrera, al tutor y a los docentes implicados en el desarrollo de la investigación. Seguidamente agradecer a familiares o a aquellas personas que se quiera. También puede incluirse en la siguiente hoja una dedicatoria personal. A modo de ejemplo el contenido podría ser:

“En primer lugar dar gracias a la Universidad Nacional de Tres de Febrero (UNTREF), a su Rector Lic. Anibal Jozami, a todo su personal docente y no docente. Por promover un espacio ideal para el desarrollo de ideas y nuevos pensamientos y brindar a todos y cada uno de los alumnos, de esta casa de altos estudios, todos los recursos que esta institución dispone. Esta investigación no hubiera sido posible sin una formación académica acorde, por este motivo debo extender mi agradecimiento a los docentes de la carrera de Ingeniería de Sonido de la UNTREF, a su coordinador Ing. Alejandro Bibondo, que siendo la primera carrera de estas características del país, es muy importante contar con un cuerpo docente afín a las exigencias que este desafío propone, prestando su dedicación y vocación de enseñar. Un especial agradecimiento por la participación de esta tesis a la tutora Ing. Nombre Apellido, que supo transmitirme sus conocimientos y ayudarme a organizarme y fijarme un rumbo concreto y delineado, disponiendo desmedidamente de su tiempo. Por otra parte, quisiera hacer una mención especial al Ing. Hernan San Martin, que permitió el uso de las instalaciones de su laboratorio para poder trabajar y la disposición de todos sus recursos para que dicha investigación se realizara en tiempo y forma. Por último y no menos importante, quiero dar un afectuoso y cálido agradecimiento a mi familia...”

DEDICATORIA

Elige a quién o a qué quieres dedicárselo.

Elegir el motivo de la dedicatoria (orientativo).

ÍNDICE DE CONTENIDOS

RESUMEN	VII
ABSTRACT	VIII
1 INTRODUCCIÓN	1
1.1 FUNDAMENTACIÓN	1
1.2 OBJETIVOS	2
1.2.1 Objetivo general	2
1.2.2 Objetivo específico	3
1.3 ESTRUCTURA DE LA INVESTIGACIÓN	3
2 MARCO TEÓRICO	5
3 ESTADO DEL ARTE	6
4 DESARROLLO	7
4.1 RECOPILCIÓN DE BASES DE DATOS	7
4.1.1 Datos in-the-wild	7
4.1.2 Datos profesionales	7
4.2 DESARROLLO DE LA CADENA DE PRE PROCESAMIENTO	7
4.2.1 Diferentes configuraciones	7
4.3 EVALUACIÓN DE LOS CONJUNTOS DE DATOS	7
4.3.1 Reducción del dataset	7
4.3.2 Calidad de la grabación	7
4.3.3 Condiciones acústicas	8
4.3.4 Diferencias del habla	8
4.3.5 Métrica conjunta	8
4.4 ENTRENAMIENTO DEL MODELO DE ESTIMACIÓN DE DENSIDAD	8
4.4.1 Validación con medelo zero-shot	8
4.5 DESCRIPCIÓN DE PRUEBAS ESTADÍSTICAS	8
4.6 MODELO DE TTS ZERO-SHOT	8
5 RESULTADOS Y ANÁLISIS	9

5.1	RESULTADOS DE LAS DIFERENTES VARIANTES DE LA CADENA	9
5.2	COMPARACIÓN ENTRE CONJUNTOS DE DATOS	9
5.3	COMPARACIÓN POR ESTIMACIÓN DE DENSIDAD	9
5.4	VALIDACIÓN MODELO DE TTS ZERO-SHOT	9
6	CONCLUSIONES	10
7	LÍNEAS FUTURAS DE INVESTIGACIÓN	11

Índice de figuras

Índice de tablas

RESUMEN

Su contenido no debe superar una página. Se indicarán los objetivos del trabajo, los métodos y resultados principales. A dos espacios debajo del resumen, en la misma página, se colocarán hasta 5 palabras clave que identifican los contenidos del trabajo.

Palabras Clave:

ABSTRACT

Ídem que para castellano.

Keywords:

1. INTRODUCCIÓN

1.1. FUNDAMENTACIÓN

Los modelos de texto a habla (TTS, por sus siglas en inglés) experimentaron un avance tecnológico exponencial en los últimos años: mediante redes neuronales profundas se alcanzaron resultados de elevada calidad sonora e inteligibilidad (Tan et al., 2021). No obstante, la fuerte dependencia de estos sistemas respecto a los datos de entrenamiento dificulta la obtención de voces sintetizadas con naturalidad para la gran diversidad de hablantes. Esta dificultad es especialmente notable en regiones con escasez de conjuntos de datos extensos, como ocurre en distintas provincias de Argentina.

En este marco, se han desarrollado sistemas de TTS en español rioplatense (Ortega Riera et al., 2023) que alcanzan resultados aceptables, pero se enfrentan a la limitada cantidad de datos específicos de los diferentes dialectos de Argentina, lo cual impide lograr sistemas más robustos y naturales. Tradicionalmente, la generación de bases de datos para entrenar modelos de TTS se orienta a recopilar grandes volúmenes de grabaciones de alta calidad (realizadas en estudios profesionales) y a emplear hablantes con características específicas (por ejemplo, locutores), lo que da lugar a un corpus homogéneo en sus características acústicas y prosódicas. Este enfoque fue crucial para la convergencia de modelos basados en aprendizaje profundo, pero representa una barrera de entrada para numerosos idiomas y variedades dialectales que no disponen de recursos para producir dichos datasets.

La literatura denomina “idiomas de bajos recursos” (low-resource languages) a estos casos; dentro de ellos se incluyen dialectos específicos de una lengua, como sería el español rioplatense o las variantes propias de determinadas provincias argentinas. Para entrenar modelos de TTS en lenguajes de bajos recursos se ha explorado la utilización de datos recolectados en Internet (Cooper, 2019), conformando conjuntos heterogéneos procedentes de diversas fuentes y de calidad de audio variable. Estos corpus suelen denominarse datos salvajes (ITW, “in-the-wild” por sus siglas en inglés). Además, con el avance de la inteligencia artificial generativa, han surgido diferentes mejoras en la arquitecturas de los sistemas de TTS mas actuales (Xie et al., 2024), lo que hace que los conjuntos de datos ITW sean una fuente especialmente atractiva para capturar la gran diversidad del fenómeno del habla.

El principal problema de entrenar modelos de TTS con conjuntos ITW es la elevada variabilidad en la calidad de las grabaciones, lo que incide directamente en la capacidad de los modelos neuronales para aprender los patrones subyacentes y, en muchos casos, impide la convergencia hacia resultados satisfactorios. Para abordar esta limitación, recientemente se han propuesto cadenas de preprocesamiento que extraen, a partir de un gran conjunto de datos, subgrupos con mejor calidad de audio (Yu et al., 2024). Si bien existen distintas variantes de estas cadenas en la literatura, no se ha llevado a cabo una caracterización acústica exhaustiva de la variabilidad que generan los conjuntos resultantes tras su aplicación. La validación suele basarse en el entrenamiento de modelos TTS y en la evaluación de su convergencia; sin embargo, no se suele caracterizar toda la cadena mediante parámetros acústicos que permitan comparar diferentes implementaciones bajo criterios comunes, ni definir configuraciones óptimas según objetivos distintos (por ejemplo, maximizar la calidad del audio frente a maximizar la cantidad de horas del corpus). El impacto de la calidad de los datos en el entrenamiento de modelos de TTS a sido profundamente estudiado (Ayllón et al., 2019), pero no se ha analizado las diferencias entre los dataset ITW y los dataset profesionales mediante un análisis objetivo.

La investigación propuesta en este trabajo tiene como objetivo determinar si es posible cuantificar la eficacia de estas cadenas de procesamiento mediante parámetros acústicos. Este tipo de análisis no solo facilita la iteración y la optimización de los procesos de filtrado de datos, sino que también abre la posibilidad de desarrollar con mayor facilidad bases de datos para lenguajes de bajos recursos, contribuyendo así a disponer de sistemas TTS de mayor calidad para una amplia variedad de idiomas y acentos locales.

1.2. OBJETIVOS

1.2.1. Objetivo general

El objetivo de la investigación es evaluar con parámetros objetivos y subjetivos, el impacto de cadenas de procesamiento de conjuntos de datos *in-the-wild* para el entrenamiento de modelos de texto a voz basados en redes neuronales profundas.

1.2.2. Objetivo específico

Los objetivos específicos son:

- Crear un dataset *in-the-wild* en español de Argentina. Recopilar datasets de voces profesionales en español (grabaciones de alta calidad realizadas por hablantes profesionales).
- Desarrollar una cadena automática de preprocesamiento modular para la generación de conjuntos de datos de habla, y procesar el conjunto de datos ITW con la cadena bajo diferentes configuraciones operativas.
- Evaluar métricas acústicas en los distintos conjuntos de datos generados y comparar dichos resultados con los obtenidos en datasets tradicionales y determinar, según criterios acústicos, cuál de los conjuntos generados puede considerarse óptimo (comparando media y desvío de los diferentes conjuntos).
- Entrenar un modelo de estimación de distribuciones y comparar la similitud entre los diferentes conjuntos en el espacio latente. Determinar el conjunto de datos óptimo según criterios de similitud basados en estimación de densidad.
- Comparar los resultados del análisis acústico con los derivados del análisis por estimación de densidad. Analizar estadísticamente la relevancia de las diferencias observadas en los distintos parámetros.
- Validar los resultados en el contexto de clonación de voz mediante modelos TTS zero-shot.

1.3. ESTRUCTURA DE LA INVESTIGACIÓN

El trabajo propuesto se enmarca en una investigación de tipo tecnológica, donde se busca desarrollar y evaluar una herramienta de software que permita la selección automática de audios para la generación de conjuntos de datos del habla, con el objetivo de generar un dataset en español de Argentina. Este desarrollo es fundamental para el desarrollo de tecnologías del habla en Argentina y contribuir así a la soberanía tecnológica nacional. El desarrollo de esta tesis se enmarca dentro del proyecto Archivoz del grupo de investigación Intercambios Transorgánicos, radicado en el MUNTREF.

El documento presenta la siguiente organización:

En el capítulo 2 se presenta el marco teórico donde se explican los fundamentos de inteligencia artificial, centrandonos en las arquitecturas que se aplican a los modelos de TTS modernos. Luego se detallan las métricas acústicas elegidas para caracterizar.

En el capítulo 3 se hace una recapitulación de los modelos de TTS actuales y de las cadenas de procesamiento que aparecieron en los últimos años.

En el capítulo 4...

2. MARCO TEÓRICO

En el Marco Teórico debemos incorporar la bibliografía, artículos de revistas, ponencias de congresos, links de Internet o todo aquello que haya contribuido a formar el cuerpo del saber sobre el que va a basarse la investigación, incorporando los procesos y ecuaciones necesarios.

Puede ser uno varios capítulos donde se detallen los parámetros, indicadores y conceptos teóricos referentes al tema a tratar. Se recomienda no utilizar conceptos muy básicos, como definición de nivel de presión sonora, ponderación A, etc.

3. ESTADO DEL ARTE

Puede ser uno o varios capítulos que desarrollen el estado del arte del área de conocimiento donde se inserta la tesis. La profundidad del enfoque en el tratamiento de los temas debe ser adecuado para el entendimiento posterior de los resultados y discusiones de la tesis. No es necesario que sea autocontenido, es recomendable el uso amplio de referencias a trabajos previos que se encuentren en la literatura abierta sobre el tema.

4. DESARROLLO

4.1. RECOPIACIÓN DE BASES DE DATOS

4.1.1. Datos in-the-wild

Detallar el proceso de recolección de datos ITW.

4.1.2. Datos profesionales

Detallar cuales son los dataset profesiones a evaluar y porque elegí esos

4.2. DESARROLLO DE LA CADENA DE PRE PROCESAMIENTO

Describir como es la cadena de pre procesamiento

4.2.1. Diferentes configuraciones

Describir las diferentes configuraciones evaluadas

4.3. EVALUACIÓN DE LOS CONJUNTOS DE DATOS

Explicar de forma general las métricas que se van a utilizar y como determinar cual es el mejor dataset.

4.3.1. Reducción del dataset

Como medir tamaño del dataset

4.3.2. Calidad de la grabación

PESQ, STOI, SI-SDR, SNR

4.3.3. Condiciones acústicas

T30, C50, D50

4.3.4. Diferencias del habla

F0-STD, MCD

4.3.5. Métrica conjunta

Explicar como combinar toda esta info

4.4. ENTRENAMIENTO DEL MODELO DE ESTIMACIÓN DE DENSIDAD

Explicar porque necesito hacer lo del modelo de estimación de densidad

4.4.1. Validación con medelo zero-shot

Definir modelo y explicar la justificación de este experimento

4.5. DESCRIPCIÓN DE PRUEBAS ESTADÍSTICAS

Detallar un poco las validaciones estadísticas que se van a realizar

4.6. MODELO DE TTS ZERO-SHOT

Definir el modelo de TTS a usar, el porque de la selección y la descripción del último experimento

5. RESULTADOS Y ANÁLISIS

5.1. RESULTADOS DE LAS DIFERENTES VARIANTES DE LA CADENA

Comparar el análisis objetivo y determinar el mejor subset

5.2. COMPARACIÓN ENTRE CONJUNTOS DE DATOS

Comparar entre conjuntos de datos profesionales vs datos ITW (por parametros acústicos)

5.3. COMPARACIÓN POR ESTIMACIÓN DE DENSIDAD

Comparar entre conjuntos de datos profesionales vs datos ITW (por estimación de densidad)

5.4. VALIDACIÓN MODELO DE TTS ZERO-SHOT

Compara calidad de clonación respecto a la calidad o subset del audio original

6. CONCLUSIONES

En las conclusiones del Plan de Investigación, debe plantearse cómo será la exposición de los resultados y qué es lo que se espera obtener en resumen de las pruebas que se realicen.

7. LÍNEAS FUTURAS DE INVESTIGACIÓN

Este trabajo pretende contribuir a la unificación de criterios en el diseño y evaluación de cadenas de preprocesado, lo que facilitará la identificación de las configuraciones más adecuadas para distintos casos de uso. Entre las líneas futuras de investigación se destacan, en particular, el desarrollo de modelos de denoising o speech enhancement personalizados: mediante técnicas de adaptación, dichos modelos buscarían aproximar conjuntos de grabaciones diversas hacia el dominio acústico del corpus con el que fue entrenado el modelo TTS, aumentando así la robustez y la fidelidad de la síntesis en condiciones reales.

Bibliografía

- Ayllón, D., Sánchez-Hevia, H., Figueroa, C., & Lanchantin, P. (2019). Investigating the Effects of Noisy and Reverberant Speech in Text-to-Speech Systems. *Proc. Interspeech*, 1511-1515. <https://doi.org/10.21437/Interspeech.2019-3104>
- Cooper, E. (2019, enero). *Text-to-Speech Synthesis Using Found Data for Low-Resource Languages* [PhD. thesis]. Columbia University.
- Ortega Riera, P., Passano, N., Paez, D., Bach, F., Pupkin, I., Sacerdoti, E., Yommi, M., & Martín, H. (2023). Implementación y Evaluación de un Sistema de Clonación de Voz Rioplatense para Asistencia en la Comunicación Oral. *Jornadas de Acústica, Audio y Sonido*.
- Tan, X., Qin, T., Soong, F. K., & Liu, T.-Y. (2021). A Survey on Neural Speech Synthesis. *arXiv preprint arXiv:2106.15561*. <https://arxiv.org/abs/2106.15561>
- Xie, T., Rong, Y., Zhang, P., & Liu, L. (2024). Towards Controllable Speech Synthesis in the Era of Large Language Models: A Survey. *arXiv preprint arXiv:2412.06602*. <https://arxiv.org/abs/2412.06602>
- Yu, J., Chen, H., Bian, Y., Li, X., Luo, Y., Tian, J., Liu, M., Jiang, J., & Wang, S. (2024). Auto-Prep: An Automatic Preprocessing Framework for In-The-Wild Speech Data. *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1136-1140. <https://doi.org/10.1109/ICASSP48485.2024.10447759>

ANEXO I. FORMATO INTERNO

AI 1. Numeración

Las páginas serán enumeradas a partir del Índice de Contenidos, con números romanos colocados en la parte media inferior de cada página. A partir de la Introducción, todas las páginas serán enumeradas con números arábigos ubicados en la parte inferior derecha. No usar la palabra “página” antes de la numeración de las páginas.