

# รายงานผู้เชี่ยวชาญ: แนวทางและแบบจำลอง AI สำหรับการฟื้นฟูและสังเคราะห์เสียงพูดจากเทปอนาคตอีกเก่า

รายงานฉบับนี้จัดทำขึ้นเพื่อนำเสนอแนวทางเชิงวิเคราะห์และปฏิบัติในการสร้างเสียงพูดที่ "เหมือนเสียงเดิมแต่ชัดเจน" จากไฟล์เสียงที่บันทึกด้วยเทปสมัยเก่า ซึ่งมีความเสียหายทางกายภาพและสัญญาณรบกวนในระดับสูง การบรรลุเป้าหมายด้านความชัดเจนระดับสูง (High-Fidelity) จำเป็นต้องใช้กลยุทธ์แบบบูรณาการที่เรียกว่า **Hybrid Workflow** โดยการรวมการประมวลผลสัญญาณดิจิทัล (DSP) เฉพาะทางเข้ากับการสังเคราะห์เสียงขึ้นใหม่โดยใช้แบบจำลองปัญญาประดิษฐ์ (AI Re-synthesis) ขั้นสูง

## 1. การวินิจฉัยความเสียหายของเสียงจากเทปอนาคตอีกและครอบแนวคิดการฟื้นฟู

การบันทึกเสียงจากเทปแม่เหล็กทำมักประสบปัญหาความบกพร่องเฉพาะทางที่แตกต่างจากสัญญาณรบกวนทั่วไปในสภาพแวดล้อมดิจิทัล

การเข้าใจความบกพร่องเหล่านี้เป็นขั้นตอนแรกที่สำคัญในการเลือกเครื่องมือฟื้นฟูที่ถูกต้อง

### 1.1. การจำแนกชั้นของการบกพร่องของเทปแม่เหล็ก (Analog Tape Artifacts)

ความบกพร่องหลักที่พบในเทปเก่าประกอบด้วย:

#### 1.1.1. Tape Hiss และ Noise Floor

เสียงชา (Tape Hiss) เป็นเสียงรบกวนความถี่สูงพื้นฐานที่เกิดจากขนาดของอนุภาคมีเหล็กบนเทป<sup>1</sup> ในอดีตได้มีการพัฒนาระบบลดสัญญาณรบกวน เช่น Dolby Noise Reduction (NR) ซึ่งทำงานโดยใช้เทคนิค Companding (การบีบอัดไดนามิกขณะบันทึกและขยายกลับเมื่อเล่น) เพื่อลดเสียงชา<sup>1</sup> อย่างไรก็ตาม เทปเก่าที่เสื่อมสภาพหรือไม่ได้บันทึกด้วยระบบ NR ดังกล่าวจะยังคงมี Hiss ที่โดดเด่น ซึ่งปัจจุบันเครื่องมือ AI-driven เช่น iZotope RX สามารถเรียนรู้ลักษณะเฉพาะของ Hiss และลดออกได้อย่างมีประสิทธิภาพ<sup>2</sup>

#### 1.1.2. Wow และ Flutter (ความผันผวนของระดับเสียง)

ความบกพร่องนี้ถือเป็นอุปสรรคทางเทคนิคที่ร้ายแรงที่สุดในการฟื้นฟูเสียงพูดจากเทป โดยถูกกำหนดให้เป็นการmodulationความถี่ปรสิต (Parasitic Frequency Modulation) ซึ่งผู้ฟังรับรู้เป็นการแกว่งของระดับเสียง (Pitch Fluctuation)<sup>4</sup> Wow หมายถึงความผันผวนที่ช้าและยาว ในขณะที่ Flutter คือความผันผวนที่รวดเร็ว<sup>5</sup> ปัญหาเหล่านี้เกิดขึ้นจากความไม่เสถียรของความเร็วของเตอร์ ความเสียหายของเทป หรือกระบวนการตัดต่อที่ไม่เหมาะสม<sup>4</sup>

การแก้ไข Wow และ Flutter

นั้นมีความสำคัญอย่างยิ่งและต้องทำเป็นลำดับแรกก่อนการกำจัดสัญญาณรบกวนทั่วไป เนื่องจากมันบิดเบือนโครงสร้างโทนเสียงหลัก (Tonal Structure) ของเสียงพูด<sup>4</sup> เครื่องมือเฉพาะทาง เช่น Celemony Capstan ใช้เทคโนโลยี DNA Direct Note Access (คล้ายกับที่ใช้ใน Melodyne) เพื่อตรวจจับและแก้ไขความผันผวนของความเร็วโดยใช้การปรับความเร็ว (Varispeed) ซึ่งให้ผลลัพธ์ที่ปราศจากความผิดเพี้ยนจากการเปลี่ยนระดับเสียง (Pitch-shifting) หรือการยืดเวลา (Time-stretching)<sup>6</sup> iZotope RX ก็มีโมดูล Wow & Flutter ที่ออกแบบมาสำหรับการแก้ไขปัญหานี้โดยเฉพาะ<sup>5</sup>

### 1.1.3. Bandwidth Limitation และ Distortion

การบันทึกเสียงเก่ามักมีข้อจำกัดด้านความกว้างของแอบความถี่ (Low Bandwidth) ทำให้เสียงขาดความคมชัดและความละเอียดสูง (High-Resolution)<sup>7</sup> ซึ่งเป็นข้อจำกัดทางกายภาพที่ไม่สามารถแก้ไขได้ด้วยการ Denoising ทั่วไป แต่ต้องอาศัยเทคโนโลยี Audio Super Resolution ที่ขับเคลื่อนด้วย AI เพื่อเติมเต็มช่องว่างความถี่ที่ขาดหายไป<sup>8</sup>

## 1.2. ภาพรวมของแนวทางแก้ไขด้วย AI (Dual-Pathway Approach)

การฟื้นฟูเสียงที่เสียหายอย่างหนักให้มีความชัดเจนเทียบเท่าระดับสตูดิโอต้องใช้กลยุทธ์สองขั้นตอนหลักที่ทำงานร่วมกัน:

- **Pathway A: Direct Speech Enhancement (การปรับปรุงเสียงพูดโดยตรง):** มุ่งเน้นไปที่การทำความสะอาดสัญญาณรบกวน (Denoising) และการแก้ไขความบกพร่องทางกายภาพ (Artifact Correction) จากไฟล์เสียงต้นฉบับ<sup>3</sup> เพื่อปรับปรุงอัตราส่วนสัญญาณต่อสัญญาณรบกวน (SNR)
- **Pathway B: Voice Cloning and Re-synthesis (การจำลองเสียงและสังเคราะห์ใหม่):** นี่คือแนวทางที่ทรงพลังที่สุดในการสร้างเสียงที่ "คงชัดราวกับอัดในสตูดิโอ"<sup>11</sup> โดยการสกัดเอกลักษณ์เสียง

(Timbre) ที่จำเป็นออกจากสัญญาณที่อาจจะยังไม่สมบูรณ์

แล้วนำไปใช้เป็นเงื่อนไขในการสังเคราะห์คำพูดใหม่ทั้งหมดด้วยโมเดล Text-to-Speech (TTS) คุณภาพสูง

12

การดำเนินการตามกลยุทธ์นี้อย่างมีประสิทธิภาพต้องตระหนักถึงความจำเป็นของ **Hybrid Workflow**

การทำความสะอาดเชิงวิเคราะห์เสียง (Pathway A) โดยเฉพาะการแก้ไข Wow/Flutter และ Hiss เป็นต้น ถือเป็นขั้นตอน Pre-processing ที่ขาดไม่ได้ก่อนที่จะใช้ Pathway B (การสังเคราะห์ใหม่)

เนื่องจากโมเดลการโคลนเสียงขั้นสูงต้องการอินพุตที่สะอาดเพื่อป้องกันไม่ให้มีการโคลนความบิดเบือนของระดับเสียง (Pitch distortion) ซึ่งเป็นสัญญาณรบกวนที่ปะปนอยู่ในเอกสารภาษาเดิมที่ผู้พูดเข้าไปด้วย<sup>12</sup>

## 2. แนวทางที่ 1: การฟื้นฟูสัญญาณเสียงโดยตรงด้วย AI (AI-Powered Speech Enhancement)

Pathway นี้ใช้ Deep Neural Networks (DNNs) เพื่อปรับปรุงคุณภาพสัญญาณเดิม

โดยเหล่านี้แทนที่วิธีการประมวลผลสัญญาณแบบดั้งเดิมด้วยการเรียนรู้แบบ End-to-end

เพื่อลดสัญญาณรบกวนได้อย่างมีประสิทธิภาพมากขึ้น<sup>14</sup>

### 2.1. สถาปัตยกรรมแบบจำลองหลักสำหรับ Speech Enhancement (Denoising)

Deep learning ได้สร้างการเปลี่ยนแปลงครั้งใหญ่ในการลดสัญญาณรบกวนและเสียงสะท้อน<sup>10</sup>

โดยมีการเปรียบเทียบประสิทธิภาพของโมเดล State-of-the-Art (SOTA) อย่างต่อเนื่อง:

- **U-Net และ Wave-U-Net:**

สถาปัตยกรรมเหล่านี้ได้รับการพิสูจน์แล้วว่ามีประสิทธิภาพสูงในการปรับปรุงสัญญาณรบกวน (Noise Suppression)<sup>15</sup> U-Net แสดงให้เห็นการปรับปรุง SNR (Signal-to-Noise Ratio) ที่สูงมาก (เช่น +71.96% บน SpEAR dataset) ทำให้เสียงรบกวนพื้นหลังลดลงอย่างเห็นได้ชัด<sup>15</sup>

- **CMGAN (Conditional Multi-Scale Generative Adversarial Network):**

แบบจำลองนี้มีความโดดเด่นในการให้ความสำคัญกับคุณภาพการรับรู้ของเสียงมนุษย์ (Perceptual Quality) โดยบรรลุคุณภาพ PESQ (Perceptual Evaluation of Speech Quality) ที่สูงที่สุด<sup>15</sup> ซึ่งหมายความว่าแม้เสียงรบกวนจะถูกลดลง

แต่เสียงพูดที่เหลืออยู่จะฟังดูเป็นธรรมชาติและสามารถเข้าใจได้ง่าย (natural and intelligible speech) ซึ่งหมายความว่าเมื่อฟังแล้วเสียงดูเป็นธรรมชาติของเสียงพูดเดิม

ในการเลือกระหว่างโมเดลลดสัญญาณรบกวนนั้น ต้องคำนึงถึงความสมดุลระหว่างการลดสัญญาณรบกวน (Noise

Suppression) และการรักษาคุณลักษณะเฉพาะของเสียงผู้พูด (Speaker-specific feature retention) ซึ่งสามารถวัดได้ด้วยคะแนน VeriSpeak<sup>15</sup> หากเลือกโมเดลที่เน้นการลดสัญญาณรบกวนมากเกินไป (เช่น U-Net ที่ทำ SNR ได้สูงสุด) อาจทำให้เกิดความบกพร่องของเสียง (Acoustic Artifacts) หรือทำให้เสียงพูดฟังดูเป็นหุ่นยนต์ (Robotic Artifacts) ได้ ในทางกลับกัน Wave-U-Net และ CMGAN แสดงให้เห็นถึงความสามารถในการรักษาเอกลักษณ์เสียงได้ดีกว่า ทำให้เป็นทางเลือกที่สมดุลกว่าสำหรับการฟื้นฟูเสียงที่ต้องการความเหมือนเดิม<sup>15</sup>

## 2.2. เทคโนโลยี Source Separation และเครื่องมือเชิงพาณิชย์

- **Source Separation:** เทคโนโลยี AI เช่นที่ใช้โดย LALAL.AI หรือโมดูลใน iZotope RX<sup>3</sup> สามารถแยกองค์ประกอบของเสียงออกจากกันได้อย่างรวดเร็วและแม่นยำ<sup>16</sup> เทคโนโลยีนี้มีบทบาทสำคัญในการฟื้นฟูเสียงพูดโดยเฉพาะอย่างยิ่งเมื่อมีเสียงดนตรีหรือเสียงบรรยายภาคปะปนอยู่ ดังที่เห็นในกรณีศึกษาของเพลง "Now and Then" ของ The Beatles ซึ่งทีมงานใช้ AI เพื่อแยกเสียงร้องของ John Lennon ออกจากเสียงเปียโนในเทปสาธิตคุณภาพต่ำ ทำให้ได้เสียงร้องที่ใสราวกับแก้ว<sup>17</sup>
  - **เครื่องมือเชิงพาณิชย์:**
    - **iZotope RX:** เป็นมาตรฐานอุตสาหกรรมสำหรับการซ่อมแซมเสียง โดยใช้ Machine Learning ในการตรวจจับและแก้ไขปัญหาที่ซับซ้อน เช่น Hiss, Hum, Click, Wow & Flutter และ Dialogue Isolate<sup>3</sup>
    - **Adobe Enhance Speech:** เป็นเครื่องมือออนไลน์ที่ใช้ AI เพื่อยกรายดับการบันทึกเสียงพูดให้คุณภาพเทียบเท่าการอัดในสตูดิโอพอดแครสต์ระดับมืออาชีพ<sup>11</sup>
    - **Remasterify และ LALAL.AI:** เสนอบริการออนไลน์ที่เน้นการทำจัดเสียงรบกวนเฉพาะทาง (tape hiss, vinyl crackle, room hum) โดยใช้ Source Separation และ AI Mastering<sup>16</sup>
- ตารางที่ 1: การวิเคราะห์เชิงเปรียบเทียบของแนวทางการฟื้นฟูสัญญาณเสียงโดยตรง (Pathway A)

เครื่องมือ/วิธีการ	หน้าที่หลัก	ข้อได้เปรียบทางเทคนิค	ความเกี่ยวข้องกับเสียงรบกวนแบบ	เกณฑ์วัดผลลัพธ์ที่สำคัญ
iZotope RX / Capstan	การแก้ไข DSP เนพาะทาง	การแก้ไขความผันผวนของระดับเสียง (Wow, Flutter)	การประมวลผลเบื้องต้นที่จำเป็นสำหรับความบกพร่องทางภาษา	ความเสถียรของระดับเสียง / การกำจัดความผิดเพี้ยน

		อย่างแม่นยำ <sup>5</sup>	พของอนาคต	น
U-Net	การลดสัญญาณรบกวนด้วย Deep Learning	อัตราการปรับปรุงสัญญาณรบกวนสูง (SNR Improvement) <sup>15</sup>	การกำจัดเสียงพื้นหลังคงที่ (Hiss, Hum) ได้อย่างมีประสิทธิภาพ	Noise Suppression
CMGAN	การปรับปรุงคุณภาพการรับรู้	ให้คะแนน Perceptual Quality (PESQ) สูงสุด <sup>15</sup>	เหมาะสมสำหรับการทำให้เสียงพูดที่ได้ฟังเป็นธรรมชาติ	Perceptual Quality
LALAL.AI / Source Separation	การแยกองค์ประกอบความเร็วและความแม่นยำในการแยกแหล่งเสียง <sup>16</sup>	การแยกเสียงพูดออกจากเสียงดนตรีหรือเสียงรอบข้าง	Speech Isolation	

### 3. แนวทางที่ 2: การจำลองเสียง (Voice Cloning) และการสังเคราะห์ใหม่ (Re-synthesis)

หากไฟล์เสียงต้นฉบับเสียหายมากจนการทำความสะอาดโดยตรงไม่สามารถบรรลุความชัดเจนระดับสูงได้ ขั้นตอนต่อไปคือการสังเคราะห์เสียงขึ้นใหม่ (Re-synthesis) โดยใช้เอกลักษณ์เสียงเดิม

#### 3.1. หลักการของ Neural Voice Cloning

Voice Cloning สมัยใหม่เป็นการทำความทำงานของโครงข่ายประสาทเทียม (Neural Networks)

ที่เรียนรู้คุณลักษณะเฉพาะของเสียงผู้พูด (Timbre, Tone, Rhythm, Pronunciation)<sup>9</sup>

เพื่อสร้างโมเดลที่สามารถสร้างคำพูดใหม่ได้ตามข้อความที่ป้อน<sup>12</sup>

- การสกัด Speaker Embeddings: โมเดลจะแยกเอกลักษณ์เสียง (Speaker Identity) ออกจากเนื้อหาทางภาษา (Linguistic Content)<sup>13</sup> เพื่อให้สามารถรักษา "โทนเสียง" เดิมไว้ได้ขณะที่สร้าง "เนื้อหา" ใหม่
- ข้อกำหนดด้านข้อมูลและการเตรียมการ: แม้ว่าจะมีเทคโนโลยี Zero-Shot TTS ที่ต้องการตัวอย่างเสียงสั้นเพียง 3 วินาที<sup>20</sup> แต่ผู้ให้บริการส่วนใหญ่ยังคงยืนยันว่าเพื่อให้ได้ผลลัพธ์คุณภาพสูงสุดระดับมืออาชีพ (Professional Voice Cloning: PVC) จำเป็นต้องมีข้อมูลเสียงที่ สะอาด ขั้นต่ำ 30 นาที และแนะนำ 3 ชั่วโมง<sup>12</sup> อย่างไรก็ตาม บางแพลตฟอร์มสามารถสร้างโคลนเสียงที่มีคุณภาพดีโดยใช้ข้อมูลเพียง 2-60 นาที<sup>21</sup>

ความจำเป็นที่ต้องใช้เสียงที่ สะอาด นี้ ตอกย้ำว่าการ Denoising เป็นทางตัน (Pathway A) เป็นสิ่งที่หลีกเลี่ยงไม่ได้ก่อนที่จะเข้าสู่กระบวนการโคลน

### 3.2. แบบจำลอง State-of-the-Art (SOTA) สำหรับการสังเคราะห์เสียง

แบบจำลองที่สามารถสร้างเสียงใหม่ให้มีความชัดเจนเท่าสูตรดิจิทัล และรักษาเอกลักษณ์เสียงได้อย่างแม่นยำคือ:

#### 3.2.1. Codec-based Language Models (Codec LMs)

VALL-E, VALL-E 2, และ VALL-E X เป็นตัวแทนของความก้าวหน้าครั้งสำคัญใน Zero-Shot TTS<sup>22</sup> โดยเดลแลนน์ปรับใช้สถาปัตยกรรม Large Language Model (LLM) กับงาน Text-to-Speech โดยการเปลี่ยนสัญญาณเสียงต่อเนื่องให้เป็นโโทเคนแบบไม่อต่อเนื่อง (Discrete Codes) ผ่าน Neural Audio Codec (เช่น Descript Audio Codec หรือ EnCodec)<sup>23</sup>

- **ความสามารถในการสังเคราะห์:** VALL-E

สามารถสังเคราะห์เสียงพูดส่วนบุคคลด้วยคุณภาพสูงโดยใช้ตัวอย่างเสียงของผู้พูดที่ไม่เคยเห็นมาก่อนเพียง 3 วินาที (Zero-Shot)<sup>20</sup> ความสามารถของ VALL-E 2 คือการบรรลุมาตรฐาน Human Parity ในด้านความเป็นธรรมชาติและความคล้ายคลึงของผู้พูด<sup>23</sup>

- **การประยุกต์ใช้ในการพื้นฟู:** โดยเดลแลนน์เป็นคำตอบที่ชัดเจนที่สุดสำหรับผู้ที่ต้องการเสียงที่ "ชัดเจน" เนื่องจากเป็นการสร้างคลื่นเสียงใหม่ทั้งหมด (High-Fidelity)  
ซึ่งปราศจากข้อจำกัดด้านแบบดิจิทัลหรือความบิดเบือนทางกายภาพของเทปเก่า<sup>8</sup> หากไฟล์เสียงที่ป้อนเป็น Acoustic Prompt ถูกทำความสะอาดเบื้องต้นแล้ว (ตาม Phase 2 ของ Hybrid Workflow) VALL-E จะสามารถสกัดเอกลักษณ์เสียงที่บริสุทธิ์และสังเคราะห์คำพูดได้ในเสียงเดิม (Original Timbre) ด้วยความชัดเจนระดับสูงสุดได้<sup>20</sup>

#### 3.2.2. Diffusion Models และ Audio Super Resolution

โดยเดลแพร่กระจาย (Diffusion Models) เป็นแบบจำลองเชิงกำเนิด (Generative Models) ที่ใช้การเรียนรู้เพื่อย้อนกระบวนการเพิ่มสัญญาณรบกวน<sup>26</sup> ในบริบทของเสียงพูด โดยเดลแลนน์มีบทบาทสำคัญในการสังเคราะห์เสียงคุณภาพสูงและเป็นองค์ประกอบในการปรับปรุงคุณภาพเสียง

(Speech Enhancement)<sup>26</sup> Respeecher ใช้เทคโนโลยี GAN-based audio technology (ซึ่งมีรากฐานคล้าย Diffusion) เพื่อทำ Audio Super Resolution

ซึ่งเป็นกระบวนการที่เติมเต็มช่องว่างของแบบนัดวิดร์ที่ขาดหายไปในเสียงคุณภาพต่ำ<sup>8</sup>

สิ่งนี้สำคัญอย่างยิ่งในการฟื้นฟูเทปเก่าที่มีข้อจำกัดด้านความถี่สูง

ตารางที่ 2: การเปรียบเทียบโมเดลสังเคราะห์และโคลนเสียงพูด (Pathway B)

หมวดหมู่โมเดล	ตัวแทนโมเดล	เทคโนโลยีหลัก	ข้อมูล Prompt ขั้นต่ำ	ความสามารถในการ พื้นฟู
Codec Language Model (SOTA)	VALL-E / VALL-E 2	การสร้างแบบจำลองภาษาตามรหัสเสียงไม่ต่อเนื่อง (Discrete Code Conditional Language Modeling) <sup>23</sup>	3 วินาที (Zero-Shot) <sup>20</sup>	สังเคราะห์เสียงใหม่ด้วยความเป็นธรรมชาติสูง ความคล้ายคลึงผู้พูดสูง และ High-Fidelity
Diffusion Models / Super Resolution	Respeecher, AudioCraft	Generative model, การเติมเต็มช่องว่างแบบนัดวิดร์ <sup>8</sup>	Few-Shot Cloning <sup>27</sup>	เหมาะสมสำหรับการฟื้นฟูเสียงที่ติดข้อจำกัดความถี่ต่ำ (Low Bandwidth)

โดยสรุป โมเดล VALL-E (หรือสถาปัตยกรรม Codec LM ที่เทียบเคียงได้)

คือคำตอบเชิงเทคนิคขั้นสูงที่สุดสำหรับคำถามนี้ เนื่องจากความสามารถในการสร้างเสียงพูดที่สมบูรณ์แบบขึ้นใหม่ (Re-synthesis) โดยใช้เอกลักษณ์เสียงเดิมเป็นเงื่อนไข แม้จะมีข้อมูลต้นฉบับที่เสียหายหรือจำกัดก็ตาม

#### 4. กลยุทธ์การดำเนินงานสำหรับการฟื้นฟูเสียงเทปเก่า (Hybrid Workflow)

เพื่อให้ได้ผลลัพธ์ที่ "เหมือนเสียงเดิม แต่ชัดเจน" ตามความต้องการของผู้ใช้งานอย่างแท้จริง

ต้องใช้กระบวนการทำงานแบบบูรณาการ 3 ขั้นตอน ซึ่งเป็นการผสานรวม Pathway A และ Pathway B

##### 4.1. Phase 1: การประมวลผลเบื้องต้นและแก้ไขความบกพร่องทางภาษาพาร์ท

ขั้นตอนนี้มีจุดประสงค์เพื่อเตรียมไฟล์เสียงที่สะอาดที่สุดเท่าที่จะเป็นไปได้ เพื่อให้ AI สามารถสกัดเอกลักษณ์เสียง

(Timbre) ที่บริสุทธิ์ได้ใน Phase 3

1. การแปลงเป็นดิจิทัล (Digitization): แปลงเสียงจากเทปเป็นไฟล์ดิจิทัลที่ความละเอียดสูงที่สุด (เช่น 24-bit/96kHz)
2. การแก้ไข Wow และ Flutter: ใช้ Celemony Capstan หรือ iZotope RX Wow & Flutter<sup>5</sup> ในการกำจัดการแกว่งของระดับเสียง การดำเนินการนี้ต้องทำก่อนการ Denoising ทั่วไป เนื่องจากความถี่พื้นฐานของเสียงพุด (FO) ถูกบิดเบือนตั้งแต่ต้น การ Denoising ในภายหลังจะไม่สามารถแก้ไขปัญหานี้ได้<sup>4</sup>
3. การลด Hiss เป็นต้น: ลดสัญญาณรบกวนพื้นฐาน (Noise Floor) ด้วยโมดูล Denoise หรือ De-hiss ใน iZotope RX<sup>3</sup>

#### 4.2. Phase 2: การแยกองค์ประกอบเสียงและการปรับปรุงความชัดเจน

ขั้นตอนนี้เน้นการแยกเสียงพุดให้โดยเด่นและพร้อมสำหรับการสกัดข้อมูล

1. การแยกเสียงพุด (Dialogue Isolation): ใช้ AI Source Separation เพื่อแยกเสียงพุดออกจากเสียงพื้นหลัง เสียงเพลง หรือเสียงรบกวนรอบข้างที่ยังคงเหลืออยู่<sup>16</sup> ตัวอย่างที่โดยเด่นคือการที่เทคโนโลยี AI สามารถแยกเสียงร้องของ John Lennon ออกจากเครื่องดนตรีอื่น ๆ เพื่อให้ได้เสียงที่คมชัดอย่างไม่เคยมีมาก่อน<sup>17</sup>
2. การปรับปรุงคุณภาพเสียง (Perceptual Enhancement):  
ประมวลผลเสียงพุดที่แยกได้ต่อด้วยโมเดลที่เน้น Perceptual Quality เช่น CMGAN<sup>15</sup> หรือ Adobe Enhance Speech<sup>11</sup> เพื่อเพิ่มความชัดเจนโดยไม่ทำให้เสียงพุดฟังดูเป็นหุ่นยนต์ ความสำเร็จของ Phase 3 (การสังเคราะห์เสียงใหม่) ขึ้นอยู่กับความแม่นยำในการสร้างสคริปต์ (Transcription) ซึ่งต้องอาศัยการแปลงเสียงเป็นข้อความอัตโนมัติ (ASR)<sup>22</sup> การที่เสียงพุดถูกทำความสะอาดและแยกออกจากสัญญาณรบกวนใน Phase 2 จะช่วยปรับปรุงความแม่นยำของระบบ ASR (เช่น OpenAI Whisper) ได้อย่างมาก<sup>28</sup> ซึ่งเป็นสิ่งจำเป็นในการสร้างข้อความที่ถูกต้องเพื่อป้อนเข้าสู่โมเดล VALL-E

#### 4.3. Phase 3: การคลอนเสียงและสังเคราะห์ใหม่ (Voice Cloning & Re-synthesis)

1. การสกัด Speaker Embeddings: ใช้ไฟล์เสียงที่สะอาดที่สุดจาก Phase 2 (Acoustic Prompt) ซึ่งอาจสั้นเพียง 3 วินาที<sup>20</sup>

2. การสังเคราะห์เสียงใหม่ (Re-synthesis): ป้อน Prompt และข้อความ (Transcription) เข้าสู่โมเดล Codec LM เช่น VALL-E<sup>23</sup> โดยจะสร้างคลิปเสียงใหม่ที่มีเอกลักษณ์เสียงเดิม แต่มีคุณภาพเสียงระดับ High-Fidelity สมัยใหม่ ซึ่งเป็นการขัดข้อจำกัดทางภาษาของเทปเก่าอ ก้าวโดยสิ้นเชิง<sup>8</sup> การสังเคราะห์ใหม่นี้เป็นวิธีเดียวที่จะรับประกันว่าเสียงจะ "ชัดเจน" ได้ที่ติดไปไม่ได้พยายามซ้อมแซมสัญญาณต้นฉบับที่อาจเสียหายเกินเยี่ยวยา

## 5. ข้อสรุปและคำแนะนำเฉพาะทาง

### 5.1. การประเมินคุณภาพและการวัดผล

ในการประเมินความสำเร็จของการบันการพื้นที่ ควรใช้เกณฑ์วัดที่หลากหลายเพื่อให้แน่ใจว่าได้รักษาเอกลักษณ์เสียงไว้:

- คุณภาพการรับรู้ (Perceptual Quality): วัดด้วย PESQ เพื่อประเมินความเป็นธรรมชาติและความเข้าใจง่ายของเสียงพูด<sup>15</sup>
- ความคล้ายคลึงของผู้พูด (Speaker Similarity): วัดโดยใช้ Speaker Embedding Cosine similarity scores จากเครื่องมือวิเคราะห์เสียง เช่น SpeechBrain toolkit หรือ Resemblyzer<sup>22</sup> เพื่อยืนยันว่า Timbre และ Pitch ของเสียงใหม่ยังคงเหมือนเสียงเดิม<sup>13</sup>
- ความแม่นยำในการสร้างสคริปต์: วัด Word Error Rate (WER) และ Character Error Rate (CER) ของเสียงที่สังเคราะห์โดยใช้ ASR (เช่น Whisper หรือ Wav2Vec 2.0)<sup>22</sup> WER ที่ต่ำแสดงว่าเสียงที่สร้างขึ้นปราศจากความบกพร่องทางอะcoustic Aberrations<sup>29</sup>

### 5.2. ข้อจำกัดทางจริยธรรมและความเสี่ยง

เทคโนโลยี AI Voice Cloning ที่นำเสนocommunity ความสามารถในการสร้างเสียงที่สมจริงในระดับสูง (เช่น VALL-E) นำมาซึ่งความเสี่ยงด้าน Deepfake และการปลอมแปลงเสียงอย่างร้ายแรง<sup>30</sup> ผู้ใช้งานควรตระหนักร่วมกับเสียงเพียง 3 วินาทีก็อาจเพียงพอต่อการสร้าง Deepfake ที่ใช้ในการหลอกหลวงได้<sup>31</sup> เมื่อว่าการใช้งานในกรณีจะเป็นไปเพื่อการพื้นฟูส่วนบุคคล แต่ผลลัพธ์ที่ได้จากการสังเคราะห์เสียงใหม่นั้นมีศักยภาพในการถูกนำไปใช้ในทางที่ผิดได้เช่นกัน<sup>32</sup>

### 5.3. ข้อสรุปและคำแนะนำเฉพาะ

คำถามของผู้ใช้งานที่ต้องการ "สร้างเสียงให้เหมือนเสียงเดิม แต่ชัดเจน" ซึ่งให้เห็นว่าการแก้ไขสัญญาณรบกวน (Restoration) เพียงอย่างเดียวไม่เพียงพอ การเลือกใช้แบบจำลองเชิงกำเนิด (Generative Models) จึงเป็นคำตอบที่ดีที่สุด

1. **แบบจำลองหลักที่แนะนำ (The SOTA Model):** แบบจำลองในกลุ่ม Codec Language Model เช่น VALL-E และ VALL-E 2 เป็นเทคโนโลยีที่ทันสมัยที่สุดในการสังเคราะห์เสียงใหม่ (Re-synthesis) ที่มีคุณภาพสูงและรักษากลั่กษณ์เสียงได้อย่างแม่นยำจากตัวอย่างเสียงสั้น ๆ
2. **แนวทางที่แนะนำ (The Hybrid Workflow):** การดำเนินการตามกระบวนการ 3 ขั้นตอน (Phase 1: การแก้ไข Wow/Flutter และ Hiss -> Phase 2: การแยกเสียงพื้นด้วย AI Source Separation -> Phase 3: การสังเคราะห์เสียงใหม่ด้วย VALL-E) เป็นกลยุทธ์ที่สามารถบรรลุเป้าหมายได้โดยเฉพาะอย่างยิ่งการเน้นย้ำความจำเป็นในการแก้ไขความบกพร่องทางภาษาพาร์ทใน Phase 1 ก่อน เพื่อให้แน่ใจว่าข้อมูลเอกสารลักษณ์เสียงที่ป้อนให้ VALL-E มีความบริสุทธิ์<sup>4</sup>
3. **ความแตกต่างระหว่างการฟื้นฟูและการสังเคราะห์ใหม่:** การ Denoising (Pathway A) คือการซ่อมแซมไฟล์เดิม ในขณะที่การใช้ VALL-E (Pathway B) คือการสร้างไฟล์ใหม่ ที่สมบูรณ์แบบกว่า ซึ่งเป็นการก้าวข้ามข้อจำกัดของเทคโนโลยีเดิม เช่น ดังนั้น สำหรับการฟื้นฟูเสียงที่เสียหายหนัก การสังเคราะห์ใหม่ด้วย AI จึงเป็นวิธีการที่เหนือกว่าในการมอบความชัดเจนระดับสูงสุดตามที่ผู้ใช้งานต้องการ

#### ผลงานที่อ้างอิง

1. How to reduce tape hiss while recording a cassette tape with digital data as source - Music, เข้าถึงเมื่อ พฤศจิกายน 7, 2025  
<https://music.stackexchange.com/questions/41318/how-to-reduce-tape-hiss-while-recording-a-cassette-tape-with-digital-data-as-source>
2. AI tool for background hiss removal? : r/audioengineering - Reddit, เข้าถึงเมื่อ พฤศจิกายน 7, 2025  
[https://www.reddit.com/r/audioengineering/comments/1jpea3k/ai\\_tool\\_for\\_background\\_hiss\\_removal/](https://www.reddit.com/r/audioengineering/comments/1jpea3k/ai_tool_for_background_hiss_removal/)
3. Professionally repair and enhance audio with RX 11 - iZotope, เข้าถึงเมื่อ พฤศจิกายน 7, 2025  
<https://www.izotope.com/en/products/rx.html?>

4. Sonogram of a recording with “wow” distortion - ResearchGate, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 [https://www.researchgate.net/figure/Sonogram-of-a-recording-with-wow-distortion\\_fig5\\_242675173](https://www.researchgate.net/figure/Sonogram-of-a-recording-with-wow-distortion_fig5_242675173)
5. RX 11 Wow and Flutter | Correct pitch modulations - iZotope, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.izotope.com/en/products/rx/features/wow-and-flutter>
6. Capstan - Celemony, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.celemony.com/en/capstan>
7. Breathing New Life into Old Memories: How AI is Revolutionizing Media Restoration, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://forestofdeanscanningservices.co.uk/read-more/breathing-new-life-into-old-memories-how-ai-is-revolutionizing-media-restoration>
8. Audio Super Resolution Turns Low-Quality Voice Samples into High-Quality Materials - Respeecher, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.respeecher.com/blog/audio-super-resolution-turns-low-quality-voice-samples-into-high-quality-materials>
9. Audio Super Resolution Turns Low-Quality Voice Samples into High-Quality Materials | by Respeecher, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://respeecher.medium.com/audio-super-resolution-turns-low-quality-voice-samples-into-high-quality-materials-21e5ca2d089d>
10. เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.mdpi.com/1424-8220/25/3/630#:~:text=Deep%20learning%20has%20revolutionized%20speech,an%20in%20real%2Dtime%20applications.>
11. Enhance Speech from Adobe | Free AI filter for cleaning up spoken audio, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://podcast.adobe.com/en/enhance>
12. AI Voice Cloning: Clone Your Voice in Minutes - ElevenLabs, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://elevenlabs.io/voice-cloning>
13. How to implement timbre transfer technology in speech synthesis? - Tencent Cloud, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.tencentcloud.com/techpedia/120024>
14. Creating Clarity in Noisy Environments by Using Deep Learning in Hearing Aids - PMC, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://PMC.ncbi.nlm.nih.gov/articles/PMC8463126/>
15. [2506.15000] A Comparative Evaluation of Deep Learning Models for Speech Enhancement in Real-World Noisy Environments - arXiv, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://arxiv.org/abs/2506.15000>
16. How AI Is Transforming Audio Restoration in Archival Video: from BBC to Netflix - LALAL.AI, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.lalal.ai/blog/how-ai-is-transforming-audio-restoration-in-archival-video-from-bbc-to-netflix/>
17. How AI helped the Beatles revive John Lennon's voice for their last song, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://news.northeastern.edu/2023/10/31/beatles-ai-last-song-john-lennon/>

18. AI Remastering: Preserving Vintage Recordings for a New Generation, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://blog.remasterify.com/ai-remastering-preserving-vintage-recordings-for-a-new-generation/>
19. Neural Voice Cloning with a Few Samples - NIPS papers, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <http://papers.neurips.cc/paper/8206-neural-voice-cloning-with-a-few-samples.pdf>
20. Vall E - Microsoft, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.microsoft.com/en-us/research/project/vall-e-x/vall-e/>
21. How much voice data is required to create a Local Voice Clone? - Altered AI, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.altered.ai/faqs/how-much-voice-data-is-required-to-create-a-local-voice-clone/>
22. Voice Cloning: Comprehensive Survey - arXiv, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://arxiv.org/html/2505.00579v1>
23. VALL-E - Microsoft, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.microsoft.com/en-us/research/project/vall-e-x/>
24. High-Fidelity Music Vocoder using Neural Audio Codecs - arXiv, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://arxiv.org/html/2502.12759v1>
25. Spectral Codecs: Improving Non-Autoregressive Speech Synthesis with Spectrogram-Based Audio Codecs - arXiv, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://arxiv.org/html/2406.05298v2>
26. Audio Diffusion Models in Speech Synthesis - Lightrains, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://lightrains.com/blogs/comprehensive-guide-audio-diffusion-models/>
27. State of the art in Voice Cloning: A review - Marvik, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://blog.marvik.ai/2023/03/21/state-of-the-art-in-voice-cloning-a-review/>
28. Voices from the past get an AI assist - Physical Sciences Area, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://physicalsciences.lbl.gov/2025/08/29/voices-from-the-past-get-an-ai-assist/>
29. [D] Are the hyper-realistic results of Tacotron-2 and Wavenet not reproducible? - Reddit, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 [https://www.reddit.com/r/MachineLearning/comments/845uji/d\\_are\\_the\\_hyperrealistic\\_results\\_of\\_tacotron2\\_and/](https://www.reddit.com/r/MachineLearning/comments/845uji/d_are_the_hyperrealistic_results_of_tacotron2_and/)
30. A Guide to Deepfake Scams and AI Voice Spoofing - McAfee, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.mcafee.com/learn/a-guide-to-deepfake-scams-and-ai-voice-spoofing/>
31. McAfee® Deepfake Detector flags AI-generated audio within seconds, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://www.mcafee.com/ai/deepfake-detector/>
32. A.I. Voice Cloning: Do These 6 Companies Do Enough to Prevent Misuse? - Innovation at Consumer Reports, เข้าถึงเมื่อ พฤศจิกายน 7, 2025 <https://innovation.consumerreports.org/AI-Voice-Cloning-Report-.pdf>