

١٢ - كيف تبني "روبوت" حقيقي؟

الاتصال الصوتي Voice Communication

د. علاء خميس

كلية هندسة البترول - جامعة قناة السويس

جدول رقم (١): تقسيمات الصوت

Voiced Sounds		Unvoiced Sounds	
Pure	a, e, i, o, u uh, aa, ee, er, uu, ar, aw,	Fricative	s, sh, f, th
Fricative	z, zh, v, dh	Plosive (stop)	p, t, k, h
Plosive (stop)	b, d, g	Affricate	ch
Nasal	m, n, ng	Aspirate	h
Affricate	j		
Glides	R, w, l, y		

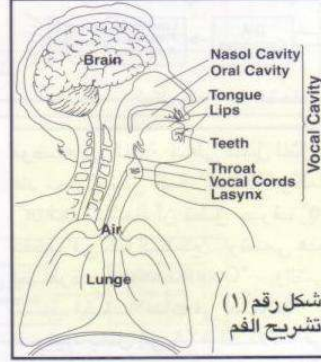
بهما الإنسان من تخمين بعض الكلمات إذا أساء نطقها المتحدث أو تخمين الغرض من الحديث أو الأمر المطلوب تنفيذه. وحالياً.. يتم تبسيط المشكلة من الفهم الكامل للغة إلى محاولة فهم بعض الأوامر الصوتية المخزنة مسبقاً في قائمة Voice Tags (أوامر مثل.. تحرك للأمام - اذهب إلى غرفة المكتب - توقف إلخ..). وهذا شبيه بما يحدث في منظومات الهواتف المحمولة والخدمات الصوتية. ويتم ذلك بتحويل الصوت إلى نص يمكن استخدامه كمدخل لبرنامج التحكم في حركة الروبوت. وإذا تطلب الأمر توفير استجابة صوتية من الروبوت للإنسان.. يتم استخدام عملية تسمى تخليق الكلام Speech Synthesis.

وتتضمن عملية الاتصال الصوتي عمليتين أساسيتين هما.. تخليق الكلام Speech Synthesis وفهم الكلام Speech Recognition. ويعتبر التخليق أسهل بكثير من فهم الكلام.. فكما هو معروف في معظم اللغات.. يمكن أن يتغير معنى الكلمة حسب موقعها في الجملة.. فإذا تخيلنا روبوت مبرمجاً على فهم اللغة الانجليزية جرت مخاطبته بالعبارة "Time Flies Like an Arrow" فسلاحظ هنا أن الكلمات الأولى "Time" و "Fly" و "Like" في هذه العبارة يصلح كل منها ليكون فعلاً للجملة. وعلى الرغم من قدرة الإنسان على فهم معنى هذه الجملة بسهولة.. إلا

الحاسب لكافة المستخدمين. فعلى سبيل المثال.. يمكن للمستخدم القيام بمهام على الكمبيوتر أثناء قيادة السيارة من خلال الصوت والحصول على رد فعل صوتي بدلاً من النظر إلى الشاشة. ومن الاتجاهات الجديدة أيضاً في مجال الاتصال الصوتي.. استخدام الصوت للتحكم في الأجهزة المنزلية كالتلفزيون والفيديو والمكيفات.. وظهر في الأسواق أجهزة «ريموت كترول» جديدة تعمل بالصوت مثل InVoca و Hands-Free Voice Accenda و Promptu.

ويستطيع Accenda مثلاً التعرف على أكثر من ٥٠ أمر لتشغيل أو إيقاف التلفزيون أو الفيديو بالإضافة إلى وظائف أخرى مثل تغيير محطة التلفزيون أو رفع الصوت أو خفضه أو بدء التسجيل على الفيديو إلخ.. ويجدر بالذكر.. أن المشكلة الأساسية التي تواجه هذه الأجهزة هي كيفية جعل جهاز الريموت كترول قادراً على التفريق بين صوت المتحدث والصوت الصادر من التلفزيون.. وهو ما يتطلب تصميم دوائر خاصة لإزالة الضوضاء.

وفي المنظومات الروبوتية.. يمثل الاتصال الصوتي تحدياً كبيراً على الرغم من الجهود التي يبذلها الباحثون لحل مشكلة فهم اللغات الحية لتحقيق الاتصال الطبيعي بين الروبوت والإنسان. وحتى إذا تم حل مشكلة الفهم الكامل للكلام.. فسوف تستمر مشكلتا التخمين والحس اللتان يمكن



التطبيقات الصناعية. وفي الحاسبات الشخصية.. يمكن أيضاً استخدام برنامج مثل Festival لتحويل النصوص إلى صوت.. في حين يمكن استخدام C-voice لتحويل الصوت إلى نص.

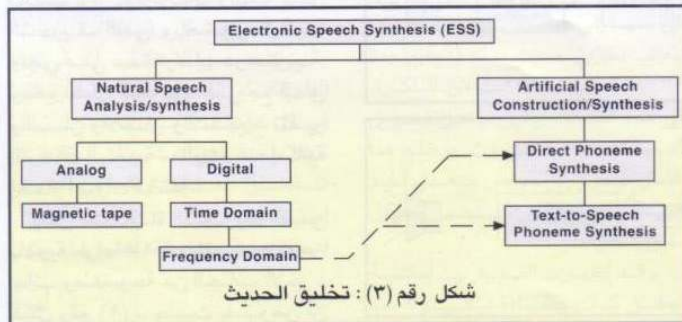
وتحاول شركة «ميكروسوفت» الآن تحسين إمكانيات الاتصال الصوتي في إصدار النوافذ الجديد والذي يحمل اسماً تجريبياً "Longhorn" وذلك من خلال وصلات اتصال طبيعية Natural User Interfaces معتمدة على الصوت. كما تحاول كثير من الشركات التي توفر خدمات البحث على الانترنت مثل «جوجل» و«ياهو» توفير خدمة البحث بالصوت بدلاً من الاعتماد على كتابة النص.. مما يسمح بالتفاعل مع محتويات شبكة الانترنت بطريقة طبيعية عن طريق إدماج الكلام مع الأنماط الأخرى من المدخلات والمخرجات. ويذكر.. أن Op-era browser هي أول أداة تصفح

تسمح للمستخدمين بتصفح الشبكة ووضع نماذج عملها بالصوت من خلال التحدث مع أجهزة الكمبيوتر الخاصة بهم.. بالإضافة إلى قيام محتويات المواقع على الشبكة بإعادة القراءة عليهم. وتعتمد هذه الأداة على استخدام تقنية الكلام من شركة IBM التي يطلق عليها Embedded Via-Voice. وبالإضافة إلى مساعدة ذوي الاحتياجات الخاصة.. فإن مثل هذه الأدوات تسهل عملية التفاعل مع

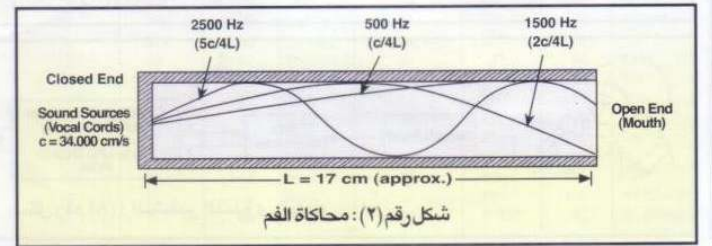
يعتبر الاتصال الصوتي من أكثر وسائل الاتصال فاعلية في حياة الإنسان.. حيث لا تقتصر لغة الحديث على تبادل معلومات بين البشر عن أحداث الحاضر أو الماضي أو المستقبل.. ولكن اللغة المنطوقة تتضمن الكثير من المعلومات الأخرى.. مثل عمر المتحدث وجنسه بالإضافة إلى حالته النفسية.

يحاول كثير من مصممي المنظومات الحاسوبية والروبوتية الآن توفير وسائل الاتصال الصوتي لجذب المستخدمين وبصفة خاصة ذوي الاحتياجات الخاصة والمعوقين مثل المكفوفين وفقادي القدرة على التحكم في أدوات الاتصال التقليدية مثل الماوس ولوحة المفاتيح وذلك بسبب اضطرابات حركية كمرضى الباركنسون Parkinson أو الشلل الرعاش حيث يعاني المريض من أعراض ببطء الحركة بالإضافة إلى التصلب أو التخشب الذي ينتج عنه فقدان الاتزان والسيطرة على الحركة. لذا.. يعتبر الاتصال الصوتي وسيلة فعالة لحل هذه المشاكل بالإضافة إلى ما يوفره من سهولة في الاتصال وإعطاء الأوامر. يستخدم الاتصال الصوتي أيضاً في بعض ألعاب الحاسب المصممة لمساعدة الأطفال المصابين بمشكلات لغوية.

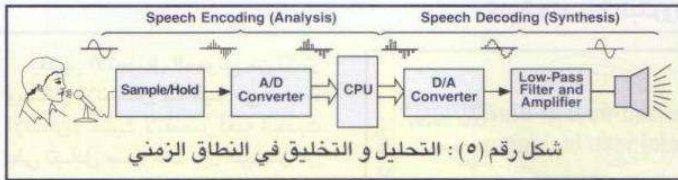
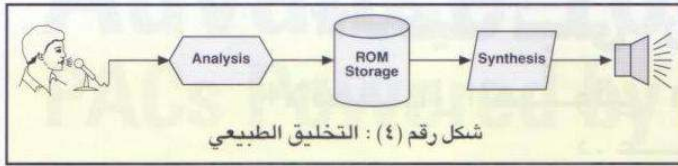
وفي مجال الاتصال الصوتي بالحاسب.. يعتبر برنامج ViaVoice الذي أنتجته شركة IBM من أكفأ البرامج التي توفر إمكانية التفاعل مع الحاسب باستخدام الأوامر المنطوقة.. ويوجد منه إصدار باللغة العربية. وقد قام باحثو المعهد الوطني الأميركي للمعايير والتكنولوجيا مع باحثين من مراكز مماثلة من فرنسا وألمانيا واليابان وبريطانيا.. بتطوير لغة برمجة جديدة أطلق عليها اسم "أيزو ١٨٦٢٩" لها القدرة على الاتصال الصوتي مع المشغل البشري لاستخدامها في



شكل رقم (٣): تخليق الحديث



شكل رقم (٢): محاكاة الفم



الصوت بتغيير التردد... مثلما يحدث بزيادة التردد عند نهاية الجملة الاستفهامية. ويسمى قالب Pattern الذي يتم فيه تغير الدرجة بالترنيم In-tonation. ويتم ضبط التوقيت أو إيقاع الجملة بإدراج وقفات Pauses بين تباينات النطق. فعلى سبيل المثال.. يساعد ضبط الإيقاع على تغير نطق كلمة "Approximate" في الجملتين "The Approximate Value is .." و "To Approximate the Value, you must..". حيث تمثل في الجملة الأولى صفة و في الثانية فعلاً. وفي حالة عدم استخدام التوقيت.. نجد أن الصوت المتولد اصطناعياً يكون غير طبيعي.

موضعه في الكلمة. فعلى سبيل المثال.. عند محاولة نطق كلمة "Robot" وكلمة "Locker" نلاحظ أن نطق حرف "o" يختلف في كلتا الكلمتين.. وتسمى هذه الظاهرة "Coarticulation" والتي تسبب اختلافات في أصوات "٤٠ قُونيم" يصل إلى ١٢٨ اختلافًا في اللغة الانجليزية. تسمى هذه الاختلافات "تباينات النطق" أو Allo-phones والتي يتم استخدامها في عملية تخليق الكلام بدلاً من «الفونيمات» للتغلب على مشكلة اختلاف نطق كل «قُونيم» حسب موضعه في الكلمة. عند تكوين جملة من مجموعة من الكلمات.. تظهر مشكلة أخرى تتمثل في ضرورة اختلاف نطق الكلمة أو نبرة الصوت حسب نوع الجملة من حيث كونها استفهامية أو خبرية. فمثلاً.. يجب أن يتغير نطق كلمة "Right" في حالة استخدامها للإستفهام مثل "Right?" أو للتعبير عن الدهشة مثل "Right!". ويتم حل هذه المشكلة بالتحكم في نطق الكلمة الموجودة في نهاية الجملة بتغيير درجة الصوت Pitch والتوقيت Timing أو الإيقاع Rhythm. يمكن تغيير درجة

الاحبال الصوتية هي الجانب المغلق بينما يمثل الجانب المفتوح الشفتين. في هذه الحالة.. تحدث حالة الرنين عند الترددات الفردية التالي $1/4(c/L)$, $3/4(c/L)$, $5/4(c/L)$, $7/4(c/L)$ حيث c سرعة الصوت (٣٤٠٠٠ سم/ث) و L طول الأنبوبة (١٧ سم).. وبالتالي تكون ترددات الرنين هي ٥٠٠ و ١٥٠٠ و ٢٥٠٠ و ٣٥٠٠ على التوالي. وبتقريب مثال الأنبوبة بحالة التجويف الفمي.. نجد أن هذه الترددات هي الترددات الرنينية التي تجعلك تستطيع إصدار الصوت "aaaah".

يمكن تقسيم الاصوات التي يصدرها الإنسان إلى قسمين أساسيين هما.. أصوات مصدرها اهتزاز الأحبال الصوتية تسمى Voiced Sounds مثل أصوات الحروف المتحركة في اللغة الانجليزية (a, e, i, o, u).. وأصوات تنتج من جعل الأحبال الصوتية مفتوحة مع دفع الهواء خلال التجويف الفمي أو الانفي وتسمى Unvoiced Sound مثل أصوات الحروف الساكنة (s, f, p, t).. ويتم تقسيم هذين النوعين الى أنواع أخرى كما هو مبين بالجدول رقم (١).

تعتمد عملية تخليق الكلام Speech Synthesis على تجميع عناصر صوتية تخليقية يطلق عليها المصوتات «الفونيمات» Phonemes.. وهي أبسط وحدات اللغة التي عادة ما تتكون من حرف أو حرفين _ جدول رقم (١)- في حالة اللغة الانجليزية والتي تحتوي على ٢٦ حرفاً و ٤٠ «قُونيم». وتسمى عملية ربط «الفونيمات» لتكوين كلمة ما «سلسلة» أو Phoneme Concatenation وهي عملية غير بسيطة بسبب اختلاف نطق كل «قُونيم» حسب

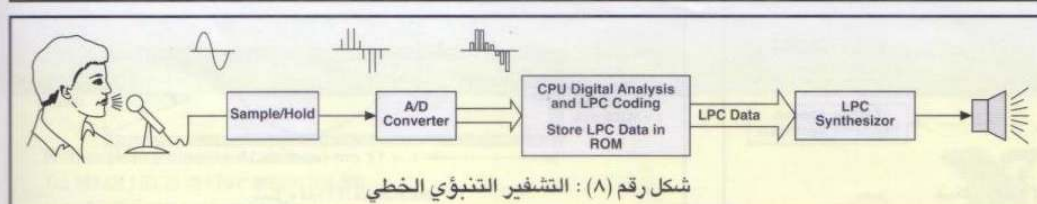
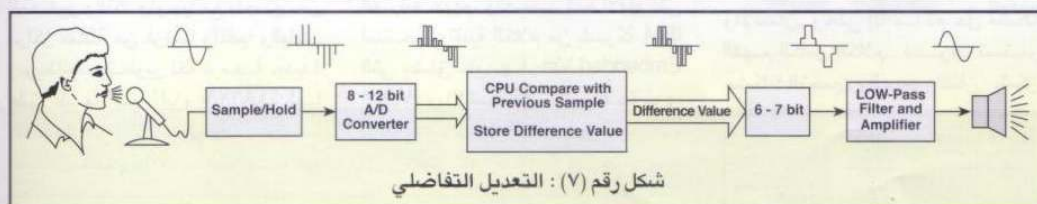
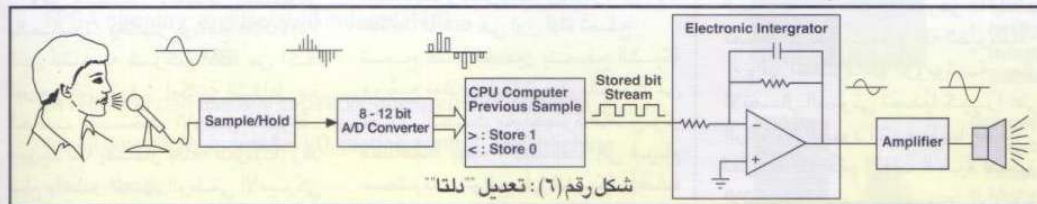
أن منظومة فهم الكلام المزود لها الروبوت يمكن أن تفهم هذه العبارة بثلاث طرق مختلفة هي "الوقت يطير مثل السهم" و "ذبابات الوقت تحب السهم" و "وقت للذبابات مثل السهم". وهناك مشكلة أخرى.. تتمثل في قدرة الروبوت على الفهم الدقيق للأوامر التي يصدرها المشغل البشري بصورة مباشرة. فمثلاً.. ثمة اختلاف بين أمر "اغسل الزجاجات قبل تعبئتها" وأمر "ادر جهاز التكيف قبل تعبئة الزجاجات" .. لا تستطيع المنظومة الحاسوبية للروبوت أن تدركه. فكلمة "قبل" في الأمر الأول.. تعني أن غسل الزجاجات يجب أن يُجْرز وينتهي قبل البدء في عملية التعبئة. أما "قبل" في الأمر الثاني.. فتعني أن تشغيل جهاز التكيف يسبق زمنياً تعبئة الزجاجات.. ولكنه سيظل مستمراً خلال عملية التعبئة. في هذه الدراسة.. يتم شرح تقنيات تخليق الكلام على أن يتم شرح عملية فهم الكلام في مقال تال بمشيئة الله.

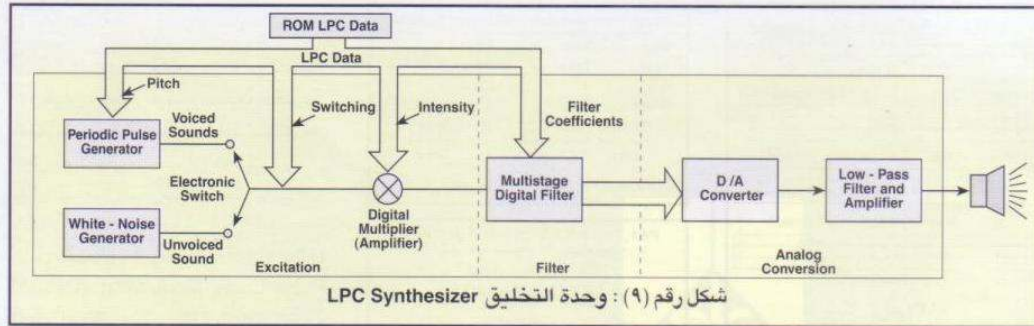
مفاهيم أساسية

قبل البدء في شرح كيفية عمل منظومات تخليق الكلام الالكترونية.. لابد من فهم الآلية التي تتم بها هذه العملية في الإنسان.. حيث تعتبر المنظومات الاصطناعية محاكاة غير كاملة لآلية الكلام التي خلقها الله عز وجل في الإنسان. يظهر الشكل رقم (١) تشريح فم الإنسان حيث يتكون من ثلاث وحدات أساسية.. هي اللسان والحنجرة والتجويف الفمي.

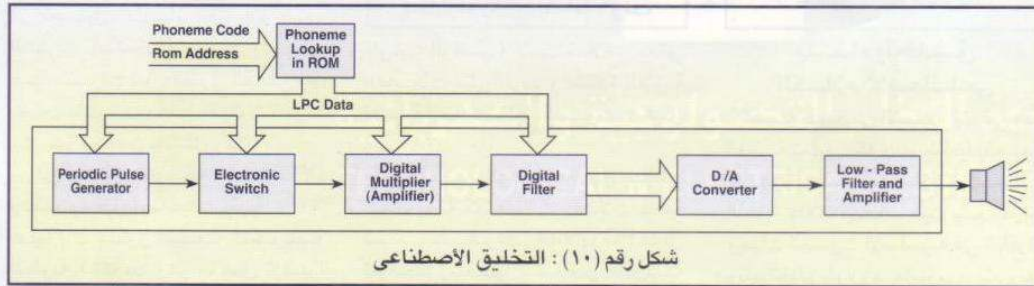
تحيل الآن أنك في عيادة طبيب حيث يطلب منك الطبيب أن تنطق كلمة "aaaah". لنطق هذه الكلمة أو لإحداث هذا الصوت.. تبدأ الرتتان في توفير القدرة اللازمة للمنظومة بدفع الهواء خلال الحنجرة إلى التجويف الفمي.. ويتولد الصوت بواسطة الأحبال الصوتية الموجودة في الحنجرة خلف تفاحة آدم والتي تتكون من طبقات جلدية تهتز عند مرور الهواء بها مما يسبب تولد عدة ترددات رنينية داخل التجويف الفمي. ويحتوى كل تردد رنيني على عدة ترددات هرمونية. وبتغير شكل التجويف الفمي مع الحلق واللسان والأسنان والشففتين.. تتغير الترددات الرنينية مما ينتج عنه إمكانية إصدار أصوات مختلفة.

يمكن محاكاة التجويف الفمي بأنبوبة طولها ١٧ سم تقريباً مغلقة من جانب ومفتوحة من الجانب الآخر - شكل رقم (٢) - حيث يفترض أن





شكل رقم (٩) : وحدة التخليق LPC Synthesizer



شكل رقم (١٠) : التخليق الاصطناعي

استخدام معدل أخذ عينات يساوى على الأقل ضعف أعلى تردد للمركبات التي تشكل الموجة الصوتية. وقد وجد.. أن مركبات الكلام البشرى عادة ما تقع بين ٣٠ - ٣٠٠٠ هرتز. لذا.. يمكن تخليق كلام بكفاءة مقبولة عند معدل أخذ عينات ٦٠٠٠ عينة/ث.

وكما ذكرنا.. فإن المعامل الآخر الذى يحدد كفاءة الكلام المخلوق هو دقة المحول التناظرى/الرقمى.. حيث يمكن الحصول على صوت مقبول فى حالة استخدام محول ٨ بت. لذا.. يكون معدل التحويل ٤٨٠٠٠ بت/ث فى حالة ضبط معدل أخذ العينات ليكون ٦٠٠٠ عينة/ث. نتيجة لذلك.. يتطلب تخزين ١٠ ثوان من الحديث.. سعة تخزينية مقدارها ٦٠٠٠٠ بايت (تقريباً ٦٤ ك بايت). ويمكن زيادة كفاءة الصوت باستخدام محول ١٢ بت.. ولكن فى هذه الحالة يجب توفير ٩٦ ك بت من الذاكرة لتخزين ١٠ ثوان.

فى طريقة تعديل «دلنا».. يتم تخزين شفرة واحدة لكل عينة بدلاً من ٨ أو ١٢ كما هو الحال فى الطريقة الأولى.. ويتم مقارنة قيمة كل عينة جديدة بالقيمة القديمة. فإذا كانت القيمة الجديدة أكبر من آخر قيمة.. يتم تخزين القيمة "١" وفى حالة العكس يتم تخزين القيمة "٠" - شكل رقم (٦).

وبالأخذ فى الاعتبار.. أنه فى هذه الطريقة يتم تخزين قيمة واحدة فقط فى كل عينة.. نجد أنه يجب زيادة معدل أخذ العينات لضمان التقاط كل تفاصيل الحديث. ولأنه من الشائع استخدام معدل ٣٢٠٠٠ عينة/ث.. فإن معدل التحويل يكون

مرحلة التحليل.. ويتم استخدام محولات رقمية/تناظرية (A/D) لإتمام عملية إعادة التخليق كما هو الحال فى منظومات الهواتف - شكل رقم (٥). تتطلب هذه التقنية قدرة تخزينية عالية لتخزين المفردات اللغوية اللازمة لتشكيل الكلمات. لذا.. ابتكرت طرق عديدة لتقليل الذاكرة المطلوبة لتخزين الصوت المشفر مثل.. طريقة تعديل الشفرة النبضية البسيطة Simple Pulse-code Modulation (PCM) تعديل دلتا Delta Modulation.. والتعديل التفاضلى Differential Pulse-code Modulation (DPCM).

فى الطريقة الأولى.. يتم التحكم فى كفاءة الصوت المشفر باستخدام معدل أخذ العينات Sampling Rate ودقة المحول التناظرى/الرقمى Resolution. وتزداد كفاءة الصوت المخلوق بزيادة معدل أخذ العينات.. ولكن على حساب زيادة سعة الذاكرة المطلوبة لإتمام هذه العملية. وكما هو معروف.. يقع الصوت البشرى فى مدى ترددى كبير يتراوح بين ١٠ - ١٥٠٠٠ هرتز. ولضمان تخزين الصوت بصورة كاملة.. يجب استخدام محول تناظرى/رقمى بمعدل أخذ عينات حوالى ٣٠٠٠٠ عينة/ث. يتم تحويل كل منها إلى شفرة ثنائية ٨ بت. لذا يتطلب تخزين حديث مدته ثانية واحدة ٨ × ٣٠٠٠٠ أى ٢٤٠٠٠٠ بت من الذاكرة. فى هذه الحالة.. يكون معدل تحويل البيانات ٢٤٠٠٠ بت/ث. يوضح هذا المثال ضرورة إنقاص معدل أخذ العينات للحصول على معدل تحويل مقبول عملياً. وقد أظهرت التجارب.. إمكانية تخليق كلام بكفاءة مقبولة عند

الذاكرة المستخدمة. هناك طرق مختلفة لتشفير الكلام.. مثل طريقة التحليل والتخليق فى النطاق الزمنى Time-Domain Analysis/Synthesis. التحليل والتخليق فى النطاق الترددى Frequency-Domain Analysis/Synthesis. التحليل والتخليق فى النطاق الزمنى: فى هذه الطريقة.. يتم تحويل التغير فى مقدار موجة الصوت التناظرية Amplitude فى شفرة رقمية باستخدام محولات تناظرية/رقمية (D/A) فى

تخزين مقاطع صوتية لأصوات بشرية مسجلة يمكن إعادة استدعائها بواسطة برنامج حاسوبى لتخليق الكلام. يمكن أن تتم عملية التخزين فى صورة تناظرية على شرائط مغناطيسية أو فى صورة رقمية فى ذاكرة الحاسب. وتنتج عملية التخزين وإعادة العرض التناظرى صوتاً شبه طبيعى ولكنها تعتبر غير عملية بسبب اقتصارها على استخدام رسائل مسجلة متسلسلة.. بالإضافة إلى استحالة تسجيل كل الكلمات والجمل التي تتطلبها عملية تخليق الكلام فى المنظومات الروبوتية.. ولكنها تعتبر أفضل اختيار عندما يكون المطلوب إنتاج لغة تخاطب محدودة كما هو الحال فى بعض المعدات أو الخدمات الهاتفية أو السيارات أو فى المطارات. تتضمن طريقة التخزين وإعادة العرض الرقمى مرحلتين أساسيتين هما.. التحليل Analysis والتخليق Synthesis - شكل رقم (٤).

فى مرحلة التحليل.. يتم تحليل الكلام البشرى وتشفيره فى صورة رقمية وتخزينه. وفى مرحلة التخليق يتم استرجاع الكلام الرقمى من الذاكرة وتحويله إلى صورة رقمية لإعادة تخليق الكلام الأصيل. توفر هذه التقنية آلية أكثر مرونة فى تخليق الكلام.. حيث لا تقتيد باستخدام كلمات أو رسائل متسلسلة أو تتابعية كما هو الحال فى الطريقة التناظرية.. بما يتيح استرجاع أى كلمة مختزنة فى الذاكرة وتتوقف عدد الكلمات على سعة

جدول رقم (٢) : خريطة الفونيمات وتباينات النطق

Hexadecimal Phoneme Code	Phoneme Symbol	Duration (ms)	Example Word	Hexadecimal Phoneme Code	Phoneme Symbol	Duration (ms)	Example Word
00	EH3	59	Jacket	20	A	185	Day
01	EH2	71	Enlist	21	AY	65	Day
02	EH1	121	Heavy	22	Y1	80	Yard
03	PA0	47	No Sound	23	UH3	47	Mission
04	DT	47	Butter	24	AH	250	Mop
05	A2	71	Made	25	P	103	Past
06	A1	103	Made	26	O	185	Cold
07	ZH	90	Azure	27	I	185	Pin
08	AH2	71	Honest	28	U	185	Move
09	I3	55	Inhibit	29	Y	103	Any
0A	I2	80	Inhibit	2A	T	71	Tap
0B	I1	121	Inhibit	2B	R	90	Red
0C	M	103	Mat	2C	E	185	Meet
0D	N	80	Sun	2D	W	80	Win
0E	B	71	Bag	2E	AE	185	Dad
0F	V	71	Van	2F	AE1	103	After
10	CH*	71	Chip	30	AW2	90	Salty
11	SH	121	Shop	31	UH2	71	About
12	Z	71	Zoo	32	UH1	103	Uncle
13	AW1	146	Lawful	33	UH	185	Cup
14	NG	121	Thing	34	O2	80	For
15	AH1	146	Father	35	O1	121	Aboard
16	001	103	Logking	36	IU	59	You
17	00	185	Book	37	U1	90	You
18	L	103	Land	38	THV	80	The
19	K	80	Trick	39	TH	71	Thin
1A	J*	47	Judge	3A	ER	146	Bird
1B	H	71	Hello	3B	EH	185	Get
1C	G	71	Get	3C	E1	121	Be
1D	F	103	Fast	3D	AW	250	Call
1E	D	55	Paid	3E	PA1	185	No Sound
1F	S	90	Pass	3F	STOP	47	No Sound



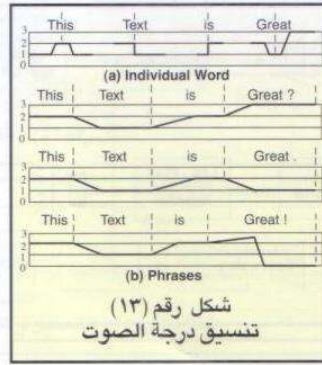
٣٢٠٠٠ بت/ث. ولتخزين حديث مدته ١٠ ثوان.. يجب توفير سعة تخزين مقدارها ٣٩ ك بت. بالمقارنة مع الطريقة الأولى.. نجد أن تعديل "دلنا" يقلل السعة التخزينية المطلوبة بواقع ٤٠٪.

في التعديل التفاضلي Differential Pulse-Code Modulation (DPCM) .. يتم استخدام نفس الفكرة المستخدمة في تعديل "دلنا" مع زيادة عدد الشفرات المستخدمة من بت واحد إلى عدة شفرات تمثل الفرق الفعلي بين عيتين متتابعين وغالباً ما يكون ٦ أو ٧ بت. كما هو مبين بالشكل رقم (٧).

في حالة منظومة DPCM-7 bit يستخدم معدل أخذ عينات بقيمة ٦٠٠٠ عينة/ث وتكون السعة التخزينية المطلوبة ٥١ ك بت وهو ما يمثل توفيراً بنسبة ١٢٪ بالمقارنة بالطريقة الأولى (PCM-8 bit) و ٢٥٪ بالمقارنة بالطريقة الثانية (DPCM-6 bit).

التحليل والتخليق في النطاق الترددي:

في هذه الطريقة.. يتم تشفير الطيف الترددي للموجة التناظرية بدلاً من تغير المقدار كما هو الحال في الطريقة الأولى.. ولا يتم تشفير الموجة التناظرية بشكل مباشر باستخدام محول تناظري/رقمي.. ولكن يتم استخدام نموذج رياضي لطيف التردد للتحكم في النموذج الإلكتروني للمنطقة الفموية للانسان Vocal Tract. في مرحلة التحليل.. يتم تحليل الخواص الترددية للصوت البشري للحصول على سلسلة من المعاملات الرياضية يتم تخزينها لاستخدامها في عملية التحكم في وحدة التخليق الإلكترونية. تحاكي هذه الوحدة المنطق الفموية للإنسان باستخدام مولدات تردد ومرشحات يتم ضبطها بواسطة المعاملات المستخلصة من مرحلة التحليل. لإتمام عملية تحليل وتخليق الكلام في النطاق الترددي.. يمكن استخدام طريقة التشفير التنبؤي الخطي Linear Pre-dictive Coding (LPC) المتميز به من القدرة على إنتاج صوت شبه طبيعي. وقد تم إقتراح هذه الطريقة بواسطة شركة "Texas Instru-ments" في لعبة تعليمية تسمى "Speak & Spell". تعتمد طريقة LPC في عملها على تشفير موجة الصوت بواسطة محول تناظري/رقمي باستخدام طريقة تعديل الشفرة النبضية البسيطة PCM المشروحة سابقاً. يتم بعد ذلك تحليل



بناء وتخليق الكلام الاصطناعي

كما يفهم من الاسم.. يتم في هذه العملية تخليق الكلام اصطناعياً بتجميع «فونيمات» أو «تباينات النطق» "Allophones". حيث يتم تخزين وحدات الصوت الأساسية في الذاكرة. وباستخدام خوارزم حاسوبي.. يتم ربط هذه الوحدات لتكوين كلمات يتم ربطها لتكوين جمل.. مع مراعاة خواص الصوت من ترنيم وإيقاع كما سبق شرحه. لإتمام هذه العملية.. يتم تشفير «الفونيمات» وتباينات النطق باستخدام طريقة تعديل الشفرة النبضية البسيطة LPC المشروحة سابقاً. ثم يتم تخزين هذه الشفرات في جدول بحث في وحدة ذاكرة ROM. يوضح الشكل رقم (١٠).. المكونات الأساسية لوحدة تخليق كلام اصطناعي تتكون من ذاكرة بحث و مصدر إثارة ومرشح متعدد المراحل. عند إدخال شفرة «فونيم» Phoneme Code على وحدة التخليق.. تقوم وحدة الذاكرة بترجمة هذه الشفرة إلى مجموعة من معاملات التعديل LPC يتم استخدامها للتحكم في مصدر الإثارة والمرشح كما هو الحال في وحدة تخليق الحديث الطبيعي المشروحة سابقاً.

يمكن استخدام هذه الطريقة في تخليق الكلام.. أو تحويل النص المكتوب إلى صوت مسموع Text-to-Speech كما هو الحال في بعض برامج الترجمة وتعلم اللغات. عند تخليق الكلام.. يتم إعداد برنامج حاسوبي لتوليد عدد من

This	PA1 = 00 111110 = 3E 1/TH = 01 111001 = 79 2/12 = 10 001010 = 8A 1/S = 01 011111 = 5F PA1 = 00 111110 = 3E
Text	2/T = 10 101010 = AA 1/EH1 = 01 000010 = 42 1/K = 01 011001 = 59 1/S = 01 011111 = 5F 1/T = 01 101010 = 6A PA1 = 00 111110 = 3E
is	1/A = 01 100111 = 67 2/Z = 10 010010 = 92 PA1 = 00 111110 = 3E
Great ?	2/G = 10 011100 = 9C 1/R = 01 011011 = 5B 3/A = 11 000000 = E0 3/T = 11 101010 = EA PA1 = 00 111110 = 3E

شكل رقم (١٥) : التشفير

Word	Pronunciation	Phoneme Symbols	
This	This	TH,I2,S	
Text	Tektst	T, EH1, K, S, T	
is	iz	I, Z	
Great	Grät	G, R, A1, T	
TH,I2,S, T,EH1,K,S,T, I,Z, G,R,A1,T			
This	Text	is	Great

شكل رقم (١١) : الشفرة الرمزية

PA1, TH, I2, S, PA1, T, EH1, K, S, T, PA1,	
This	Text
I, Z, PA1, G, R, A, T, PA1	
is	Great

شكل رقم (١٢) : إضافة الوقف

الموجة الرقمية لاستخلاص بعض المعاملات مثل التردد والشدة اللازمة لعملية إعادة تخليق الصوت - شكل رقم (٨).

يتم بعد ذلك.. وضع البيانات المستخلصة كمعاملات لمعادلات خطية تسمى شفرات LPC Codes تمثل نموذجاً رياضياً للخواص الترددية لموجة الصوت المنطوق. ويتم استخدام تلك المعاملات المستخلصة من مرحلة التحليل.. في التحكم في وحدة التخليق الإلكترونية والتي تمثل محاكاة أو نمذجة للمنطقة الفموية في الانسان - شكل رقم (٩).. وتتكون هذه الوحدة من مصدر إثارة Excitation Source .. ومرشح رقمي متعدد المراحل.. بالإضافة إلى محول رقمي/تناظري. يحتوى مصدر الإثارة على مولد نبضات دورية يقوم بمحاكاة عمل الأحبال الصوتية في الانسان عن طريق توليد نبضات دورية بترددات أصوات حروف متحركة Voiced Sound. وتحدد قيمة التردد درجة أو حدة الصوت Pitch المتولد كما هو الحال في الانسان.. حيث تتحدد درجة الصوت حسب درجة اهتزاز الأحبال الصوتية. ولتوليد أصوات حروف ساكنة Unvoiced Sounds يتم استخدام مولد الضوضاء الموضح في شكل رقم (٩). وباستخدام المفتاح الإلكتروني.. يمكن اختيار توليد أصوات حروف متحركة أو ساكنة.. ويتم تكبير الكلام المتولد بواسطة مكبر ثم يتم تمريره على وحدة الترشيح لضبط الصوت بتغيير قيم المعاملات المستخلصة من عملية التحليل لإنتاج صوت شبه طبيعي. بعد عملية الترشيح.. يتم تحويل الموجة الرقمية إلى صورة تناظرية بواسطة المحول الرقمي/التناظري.

PA1, 1/TH, 2/12, 1/S, PA1, 2/T, 1/EH1, 1/K, 1/S, 1/T,	
This	Text
PA1, 1A, 2/Z, PA1, 2/G, 1/R, 3/A, 3/T, PA1	
is	Great?

شكل رقم (١٤) : إضافة الترنيم

شفرات «الفونيمات» الواجب إرسالها إلى وحدة التخليق لتكوين كلمة ما أو جملة ما. تسمى هذه الشفرات سلسلة «الفونيمات» أو «Phoneme String». يتم تخزينها في وحدة ذاكرة RAM أو ROM كجزء من البرنامج الفرعي المسئول عن عملية التخليق. يوضح الجدول رقم (٢).. استخدام منظومة التشفير الحاسوبية السادسة عشر في بناء هذه النماذج الصوتية المنطوقة في برنامج VOTRAX لتخليق الكلام الاصطناعي.

في هذه القائمة.. يوجد ٦٤ من «الفونيمات» وتباينات النطق يتم استخدامها بواسطة برنامج فرعي لتشكيل سلسلة «الفونيمات» Phoneme String كالتالي:

١- يتم تحديد السلسلة الرمزية «الفونيمات» «Phoneme Symbol String» اللازمة لنطق كلمة معينة داخل الجملة. على سبيل المثال في الجملة "This Text is Great?".. يتم البحث في قاموس عن النطق الصحيح لكل من كلمات هذه الجملة لاستخدام هذا النطق كدليل بحث في الجدول رقم (٢).. وبالتالي يمكن تكوين الشفرة الرمزية للجملة كما هو مبين بالشكل رقم (١٤) باستخدام قاموس "Webster's New World" للغة الإنجليزية.

٢- إضافة الوقف Pauses بين المقاطع والكلمات حسب الإيقاع المطلوب. وطبقاً للجدول رقم (٢).. هناك نوعان من الوقف حسب المدة الزمنية (47-ms) PA0 و (185-ms) PA1.. يتم إضافة الوقف القصير بين مقاطع بعض الكلمات بينما يتم إدراج الوقف الطويل بين الكلمات.. يوضح الشكل رقم (١٢) الشفرة الرمزية «الفونيمات» بعد إضافة الوقف.

٣- تحديد ترنيم Intonation كل كلمة على حدة.. وترنيم الجملة بالكامل. يعتمد الترنيم على تغير درجة الصوت Pitch أو نبرة الكلام حسب موقع الكلمة ونوع الجملة. لتحديد الترنيم.. يتم تكوين مخطط بياني لتغير درجة الصوت - شكل رقم (١٣). يتم إضافة الترنيم المستنتج من الرسم البياني إلى الشفرة الرمزية «الفونيمات» كما هو مبين بالشكل رقم (١٤).

٤- تحويل السلسلة الرمزية للمصوتات إلى سلسلة شفرات «الفونيمات» باستخدام منظومة التشفير الحاسوبية السادسة عشر الموضحة بالجدول رقم (٢). يوضح الشكل رقم (١٥) سلسلة الشفرات الواجب استخدامها لنطق الجملة.

٥- استخدام سلسلة الشفرات كمدخل لوحدة التخليق الاصطناعي واستماع الجملة المنطوقة وتعديل الشفرة إذا لزم الأمر.