

## سنتز گفتار فارسی با استفاده از فرکانس گام در نرم‌افزار Flite

فاطمه نعیمی<sup>۱</sup>، دانشجوی دکتری؛ وحید قدس<sup>۲</sup>، استادیار

۱- گروه مهندسی برق - واحد سمنان - دانشگاه آزاد اسلامی - سمنان - ایران - f.naimi@stu.semnaniau.ac.ir

۲- باشگاه پژوهشگران جوان و نخبگان - واحد سمنان - دانشگاه آزاد اسلامی - سمنان - ایران - v.ghods@semnaniau.ac.ir

**چکیده:** در این پژوهش، طراحی و پیاده‌سازی سنتز کننده گفتار به زبان فارسی با استفاده از نرم‌افزار Flite ارائه شده است. بدین طریق که ابتدا توسط پایگاه داده جملات فارسی فارس‌دات<sup>۱</sup>، میانگین و انحراف معیار<sup>۲</sup> فرکانس گام و بزرگترین فرکانس گام هر واج<sup>۳</sup> واک‌دار را به دست می‌آوریم. سپس تغییرات هر واج عبارت موردنظر را از طریق محاسبه مقدار ارزش<sup>۴</sup> آن با توجه به فرمول ارائه شده پیشنهادی، در نرم‌افزار وارد می‌کنیم. ویژگی این سنتز کننده، تبدیل متن به گفتار با لهجه و تلفظ فارسی می‌باشد. در انتهای این مقاله، نتایج حاصل از اجرای الگوریتم پیشنهادی با منحنی فرکانس‌های گام به دست آمده توسط پایگاه داده جملات فارسی فارس‌دات، مقایسه شده است. همچنین، در این پژوهش چند نمونه از جملات زبان فارسی مربوط به پایگاه داده جملات فارسی فارس‌دات، با استفاده از روش پیشنهادی بر روی نرم‌افزار Flite، مورد سنتز قرار گرفته است. آزمون‌های شنیداری، برای میزان قابل فهم بودن، طبیعی بودن و خوشایند بودن جملات مذکور انجام شده است که نتایج آن‌ها برای جملات آموزش به ترتیب ۴/۴، ۴/۲ و ۴/۶ می‌باشد. همچنین برای جملات مجموعه آزمون، به ترتیب برابر ۴/۲، ۴/۱ و ۴/۳ می‌باشد.

**واژگان کلیدی:** تبدیل متن به گفتار، سنتز گفتار، فارسی، فرکانس گام، نرم‌افزار Flite.

## Farsi Speech Synthesis Using Pitch Frequency in Flite software

F. Naiemi<sup>1</sup>, PhD Student; V. Ghods<sup>2</sup>, Assistant Professor

1-Electronic Group, Semnan Branch, Islamic Azad University, Semnan, Iran, Email: f.naimi@stu.semnaniau.ac.ir

2-Young Researchers and Elite Club, Semnan Branch, Islamic Azad University, Semnan, Iran, Email: v.ghods@semnaniau.ac.ir

**Abstract:** This survey introduces a model and the implementation of a speech synthesizer in Farsi language using Flite software. In this approach, the mean and the standard deviation of pitch frequency of each voiced phoneme are first calculated by a database of Farsi sentences (Fars Dat). Then, the changes of each phoneme of the desired phrase are inserted into the software through the calculation of a value. The main feature of this synthesizer is its ability to change text to speech within Farsi pronunciation and in Farsi dialect. At the end of this paper, the results of this algorithm are compared to the changes of pitch frequencies extracted from the database of Farsi sentences. Some examples of the sentences from the database are also synthesized using our proposed method on Flite Software. The value of MOS test for understandability, naturalness and good sounding of those sentences are 4.4, 4.2, and 4.6 for the training set, respectively, and 4.2, 4.1, and 4.3 for the test set, respectively.

**Keywords:** Text to speech, speech synthesizer, Farsi (Persian), pitch frequency, Flite software.

تاریخ ارسال مقاله: ۱۳۹۶/۰۵/۱۵

تاریخ اصلاح مقاله: ۱۳۹۶/۰۸/۲۰

تاریخ پذیرش مقاله: ۱۳۹۶/۰۸/۲۶

نام نویسنده مسئول: وحید قدس

نشانی نویسنده مسئول: ایران - سمنان - دانشگاه آزاد اسلامی - واحد سمنان - باشگاه پژوهشگران جوان و نخبگان.

## ۱- مقدمه

تصور اینکه یک ماشین بتواند گفتار تولید نماید سال‌هاست که با بشر بوده است، اما تحقق وجود چنین ماشین‌هایی در طی ۵۰ سال گذشته محقق شده است. یکی از اولین کاربردهای سنتز گفتار در سال ۱۹۳۶ اتفاق افتاد. هومر دادلی [۱] در آزمایشگاه‌های بل، یک ابزار مکانیکی تولید کرد که از طریق حرکت پدال‌ها و کلیدهای مکانیکی کار می‌کرد. این ماشین می‌توانست در صورت وجود یک کاربر آموزش‌دیده، صداهایی تولید نماید (اگر داده‌های مناسبی به آن می‌رسید) که مانند گفتار به نظر آیند. این دستگاه که Voder نام داشت برای اولین بار در نمایشگاهی بین‌المللی در سال ۱۹۳۹ در نیویورک و سانفرانسیسکو به نمایش گذاشته شد [۱].

تحقق این موضوع که سیگنال‌های صوتی می‌توانند به‌عنوان مدل‌های فیلتر-منبع<sup>۵</sup> تجزیه شوند، به‌طوری‌که فاصله بین تارهای صوتی مانند یک منبع صدا و قسمت‌های مختلف دهان مانند فیلترهایی در نظر گرفته شوند، در ابزارهای الکترونیکی آنالوگ مورد استفاده قرار گرفت که می‌توانستند گفتار انسانی را تقلید نمایند. Voder توسعه‌یافته که مجدداً توسط هومر دادلی ساخته شد، می‌توانست نمونه‌ای از این ابزارها باشد [۱].

رشد سنتز الحاقی<sup>۶</sup> در دهه ۷۰ آغاز شد. در سال ۱۹۷۲ راهنمای استاندارد Unix دستورهایی برای پردازش متن به گفتار، تشکیل تجزیه و تحلیل متن، پیش‌بینی آوایی، تولید واج و سنتز شکل موج از طریق یک سخت‌افزار ویژه را ارائه نمود [۲]. ابزارهای 'n spell speak تگزاس که در اواخر دهه ۷۰ منتشر شدند، یکی از نمونه‌های اولیه تولید انبوه سنتزکننده‌های گفتار بودند [۱]. کیفیت این ابزارها در مقایسه با استانداردهای کنونی بسیار ضعیف بود؛ اما در آن زمان بسیار جالب‌توجه بودند. گفتار، اساساً با استفاده از کد کردن مبتنی بر پیش‌بینی خطی<sup>۷</sup> کدگذاری می‌شد و اغلب از کلمات و حروف جدا از هم استفاده می‌کرد. سنتزگر MITalk متعلق به دنیس کلات از بسیاری جهات مفهوم سنتز گفتار خودکار را به جهان معرفی نمود [۱]. این ابزار که بعداً به محصول DECTalk [۱] ارتقا یافت، گفتاری رباتی اما بسیار قابل‌فهم تولید می‌نمود. این ابزار یک سنتزگر فرمنت بود که تمام صنعت و هنر آن زمان را در خود جای داده بود.

در اواخر دهه ۸۰ سنتزگرهایی که فقط نرم‌افزاری بودند نیز تولید شدند، کیفیت گفتار هنوز هم شبیه صدای انسان نبود، اما گفتار می‌توانست تقریباً لحظه‌ای تولید شود. با ظهور ماشین‌های سریع‌تر و با فضای ذخیره‌سازی بیشتر، مردم شروع به ارتقای سنتزکننده‌ها با استفاده از ذخیره‌های بزرگ‌تر و متنوع‌تر برای گفتار الحاقی نمودند. یوشینوری سگیساکا [۲] در مرکز تحقیقات پیشرفته ارتباطات از راه دور ژاپن ابزار Nuu-talk را در اواخر دهه ۸۰ و اوایل دهه ۹۰ تولید نمود. این ابزار ذخیره بسیار بزرگ‌تری از واحدهای الحاقی ارائه می‌نمود. در نتیجه به‌جای یک نمونه برای هر واحد دو واجی، نمونه‌های بسیاری در اختیار بود. حتی تا سال ۱۹۹۴ زمان موردنیاز برای تولید

فایل‌های پارامتری برای یک صدای تازه در Nuu-talk (۵۰۳ جمله) تقریباً چند روز کاری یک پردازنده بود و سنتز به‌صورت لحظه‌ای امکان‌پذیر نبود.

کورزوئل [۳] پیش‌بینی کرد نسبت قیمت به عملکرد در سال ۲۰۰۵ باعث تولید سیستم‌های سنتز گفتار ارزان‌تر و قابل‌دسترس بیشتر شده و بیشتر مردم می‌توانند از این سیستم بهره‌مند شوند، همچنین از سال ۲۰۰۵ به بعد برخی از محققان با استفاده از مجموعه داده مشترک گفتار به ارزشیابی سیستم‌های سنتز گفتار پرداختند [۳]. در دسترس بودن سیستم‌های سنتز رایگان و نیمه‌رایگان مانند سیستم سنتز گفتار Flite [۴] و پروژه MBROLA [۵] هزینه‌های ورود به زمینه سنتز گفتار را کاهش می‌دهد و گروه‌های بیشتری در ارتقای این سیستم‌ها مشارکت می‌نمایند.

ازجمله تلاش‌های محققان ایرانی در مورد موضوع تبدیل متن به گفتار فارسی می‌توان موارد زیر را ذکر نمود:

۱- تبدیل حرف به صدا در زبان فارسی به کمک شبکه‌های عصبی پرسپترون چندلایه‌ای [۶] که به دلیل عدم استفاده از اعراب در نوشتار و در نتیجه مستور بودن بعضی از واژه‌ها، یک سیستم تبدیل حرف به صدا با معماری سه لایه بررسی شده است. لایه اول این سیستم قانون-گرا می‌باشد و لایه دوم از پنج شبکه عصبی پرسپترون چندلایه‌ای و یک بخش کنترل‌کننده برای تعیین دنباله واجی متناظر با حروف تشکیل شده است.

۲- تولید پارامترهای سنتز گفتار فارسی با استفاده از مدل‌های مخفی مارکوف و درخت تصمیم‌گیری [۷]. در مقاله [۷]، برای پیاده‌سازی سیستم سنتز، از مدل‌های مخفی مارکوف برای مدل‌کردن پارامترهای مربوط به واحدهای گفتاری استفاده شده است و برای تبدیل ضرایب کپسترال به سیگنال صحبت، از فیلتر MLSA استفاده شده است. برای استخراج فرکانس گام، روش خودهمبستگی اصلاح‌شده مورد استفاده قرار گرفته است. برای تولید پارامترهای سنتز گفتار توسط HMMها، از الگوریتمی استفاده شده که در آن، برای در نظر گرفتن اطلاعات بافت، علاوه بر ویژگی‌های ضرایب کپستروم و فرکانس گام، مشتق اول و دوم آن‌ها نیز، مورد استفاده قرار گرفته است [۷]. برای به‌دست آوردن مدل طول زمانی واج‌ها، مشاهدات موجود از هر تریفون را در پایگاه داده، طبق الگوریتم ویتربی با مدل HMM آن مقایسه نموده و دنباله حالات طی شده را به‌دست آورده و با میانگین‌گیری از تعداد دفعات حضور در هر حالت مدل HMM تریفون، متوسط طول زمانی حضور در هر حالت را برای هر تریفون به‌دست آمده است [۷]. زمان‌های میانگین حاصل، مدل‌های طول زمانی برای هر تریفون را تشکیل می‌دهند. در هنگام سنتز با توجه به مدل طول هر حالت از مدل HMM هر تریفون، بردار واریانس آن حالت تکرار و با استفاده از این پارامترها، دنباله ضرایب کپسترال و گام موردنیاز برای سنتز گفتار به‌دست آمده و توسط فیلتر MLSA به گفتار تبدیل

کلمه<sup>۱۱</sup> تعیین می‌گردد. در مرحله بعد با توجه به مشخصات دستوری کلمات، هر کلمه به رشته واجی مربوط به خود تبدیل می‌شود. یکی از بدیهی‌ترین راه‌ها برای تولید رشته واجی، مجموعه قواعد تبدیل حرف به صدا می‌باشد که حروف متن را به رشته‌ای از واج‌ها تبدیل می‌نماید.

## ۲-۲ مولد نوای گفتار

### ۲-۲-۱ نوای<sup>۱۲</sup>

نوا یکی از فاکتورهای اصلی برای به‌دست آوردن یک گفتار مصنوعی با کیفیت زیاد می‌باشد. مفهوم نوا، زیربوم کردن صدا و لحن گفتار است که باعث تلفظ و برداشت مفهوم‌های مختلفی از گفتار می‌شود. همچنین اطلاعاتی راجع به وضعیت روانی سخن‌گو ارائه می‌دهد. در صدای هر شخص بخش‌های صدادار<sup>۱۳</sup> و بخش‌های بی‌صدا<sup>۱۴</sup> وجود دارد. بخش‌های صدادار صدای هر شخص، دارای یک فرکانس اصلی می‌باشد که بر اساس آن تارهای صوتی به لرزه درمی‌آیند. می‌توان گفت فرکانس اصلی یکی از فاکتورهای جدانشونده صدا می‌باشد که قابل تغییر است. به‌طور ساده، فرکانس اصلی به‌وسیله اصلاح شکل‌موج قابل تغییر می‌باشد.

مفهوم مولد نوای گفتار زیربوم کردن صدا یا تغییرات فرکانس گام [۱۱-۱۳]، دیرش<sup>۱۵</sup> [۱۶-۱۴]، شدت<sup>۱۶</sup> [۱۷، ۱۸] و درنگ<sup>۱۷</sup> [۱۹] می‌باشد.

### ۲-۳ سنتز گفتار

در خصوص بلوک سنتز گفتار باید بدانیم سنتزگر گفتار مبتنی بر سه روش می‌باشد:

۱- سنتز شمرده به شمرده لغات<sup>۱۸</sup>

۲- سنتز فرمنت<sup>۱۹</sup>

۳- سنتز الحاقی<sup>۲۰</sup>

### ۲-۳-۱ سنتز شمرده به شمرده لغات

در این روش از یک مدل کامپیوتری-مکانیکی برای تولید گفتار استفاده می‌شود [۹، ۲۰]. مثل حنجره که به‌وسیله ارتعاش تارهای صوتی باعث تولید صدا و تحریکات حلقی می‌شود. زبان به عنوان بخشی از فیلتر مسیر صوتی<sup>۲۱</sup> بوده و در ایجاد اصوات نقش دارد. مسیر عبور هوا به‌وسیله زبان با فشار بر روی سقف دهان یا دندان‌ها تغییر می‌کند. لب‌ها برای ایجاد هماهنگی و کمک به فرایند تولید گفتار به زبان کمک می‌کنند. در حالت ایده‌آل سنتزکننده‌های شمرده به شمرده می‌توانند به‌وسیله ماهیچه‌های مصنوعی و کنترل آن‌ها به شکل زبان، لب و دهان شبیه‌سازی شوند.

### ۲-۳-۲ سنتز فرمنت

در دهه‌های اخیر بیشترین استفاده در متن روش‌های سنتز، روش سنتز فرمنت بوده است که بر اساس مدل فیلتر-منبع از گفتار

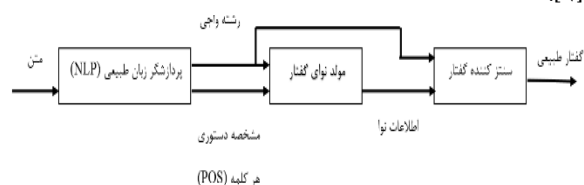
شده‌اند [۷]. برای در نظر گرفتن تأثیر پارامترهای مختلف بر نحوه تلفظ آواها، از درخت‌های تصمیم‌گیری CART، استفاده شده است. این درخت‌ها نقش تولید گام و طول زمانی واج‌ها را بر عهده دارند. برای تولید خودکار فرکانس گام از روش مطرح شده توسط فوجی ساکی [۷] استفاده شده است. در این روش برای مدل کردن شکل کلی سیگنال گام، یک جزء عمومی و برای مدل کردن تأکیدها، تعدادی اجزای محلی در نظر گرفته شده است [۷].

۳- تبدیل حرف به صدا در سیستم‌های تبدیل متن به گفتار فارسی با استفاده از درخت‌های تصمیم‌گیری CART [۲] که در آن پارامترهای استخراج‌شده از کلمات متن آموزش داده شده و سپس از درخت در تعیین آواهای تشکیل‌دهنده کلمات استفاده شده است [۲]. در ادامه مقاله، در بخش ۲ توضیحاتی در خصوص سنتزکننده گفتار ارائه شده است. در بخش ۳، تبدیل متن به گفتار در نرم‌افزار Flite بررسی شده است. بخش ۴ مدل پیشنهادی برای تولید گفتار به زبان فارسی با استفاده از نرم‌افزار Flite را معرفی می‌کند. در بخش ۵، نتایج حاصل از مدل پیشنهادی بر روی جملات فارسی ارائه شده و در پایان، بخش ۶ به نتیجه‌گیری از مباحث پرداخته است.

## ۲- سیستم سنتز گفتار

سنتز گفتار یک فناوری است که به‌وسیله آن متن به گفتار مصنوعی<sup>۸</sup> تبدیل می‌شود [۸]. در موضوع سنتز گفتار، ذخیره‌سازی تمام کلمات یک زبان غیرممکن (و اغلب بی‌فایده) است. در واقع سنتز گفتار، تولید گفتار از طریق رونویسی حروف به آوا، به‌منظور گفتن کلمات و جملات می‌باشد.

شکل ۱ روند بنیادی یک سنتز کننده گفتار را نشان می‌دهد [۹].



شکل ۱: سیستم تولید متن به گفتار فارسی [۹].

## ۲-۱ پردازش‌گر زبان طبیعی<sup>۹</sup>

در بلوک پردازش‌گر زبان طبیعی جملات ورودی به فهرستی از کلمات تبدیل می‌گردند [۱۰]. در این راستا اعداد، حروف اختصاری، علائم و نشانه‌ها نیز به معادل کامل متنی خود تبدیل می‌شوند. سپس بر روی کلمات جمله، تحلیل شکلی<sup>۱۰</sup> انجام می‌شود (در زبان فارسی و انگلیسی تطابق مستقیمی بین حروف و واج وجود ندارد، به‌عنوان مثال در زبان فارسی می‌توان به موارد زیر اشاره نمود: وجود یک حرف برای چند واج، استفاده از حروفی که خوانده نمی‌شوند، نوشته نشدن واکه‌های کوتاه، نوشته نشدن کسره اضافه) و مشخصه دستوری هر

صداهاى بی‌واک دوره تناوب ندارند، علائم گام را به‌صورت فرضی برای صداهاى بی‌واک (غیر از سکوت) در هر ۱۰ میلی‌ثانیه تعیین می‌کنیم. بدین‌صورت که یک منحنی گام به‌وسیله درون‌یابی صداهاى واک‌دار مجاور، به صدای بی‌واک (غیر از سکوت) اختصاص داده و نتیجه آن هموارسازی می‌شود؛ بنابراین در این روش سکوت از گفتار جدا می‌شود اما راهی برای تعیین علائم گام صداهاى واک‌دار از علائم گام صداهاى بی‌واک وجود ندارد [۴].

تعیین تغییرات فرکانس گام در این نرم‌افزار از طریق تعریف میانگین و انحراف معیار می‌باشد؛ که از طریق رابطه ۱ به‌دست می‌آید:

$$\text{map\_f0}(v,m,s)=(((v-mm)/ms)*s)+m \quad (1)$$

که در آن  $mm=\text{model\_mean}=170$ ،  $s=\text{stddev}$ ،  $m=\text{mean}$ ،  $v=\text{value}$  و  $ms=\text{model\_stddev}=34$  است.

در این برنامه، مقادیر  $\text{model\_mean}$  و  $\text{model\_stddev}$  ثابت می‌باشند. در این برنامه تنها می‌توانیم با تغییر میانگین و انحراف معیار که برای تمامی واج‌ها یکسان تعریف شده است، منحنی فرکانس گام تمام واج‌ها را به مقدار یک‌سان تغییر دهیم؛ اما نمی‌توانیم زیربومی هر واج را جداگانه تغییر دهیم. بنابراین صدای به‌دست‌آمده بسیار رباتیک می‌باشد و فقط با تغییر عدد مربوط به میانگین و انحراف معیار که برای کلیه واج‌ها می‌باشد، تغییرات یک‌سان به‌کار می‌رود.

واحد پایگاه داده در این سنتز کننده، دوواجی‌ها با حجم ۲/۱ مگابایت می‌باشند. در این پایگاه داده آدرس مربوط به هر یک از دوواجی‌ها مقابل آن نوشته شده است و برای هر یک پارامترهای مربوط به ضرایب LPC و سیگنال باقی‌مانده<sup>۲۴</sup> به‌دست می‌آید. با توجه به نرخ نمونه‌برداری، مرتبه ضرایب LPC در این نرم‌افزار ۱۰ می‌باشد، همچنین فرکانس نمونه‌برداری ۸ کیلوهرتز و نمونه‌ها ۱۶ بیتی هستند. Flite شامل ۴ بخش اصلی می‌باشد [۴]:

۱. کتابخانه Flite: شامل کدهای مرکزی است. این قسمت از سیستم CST نامیده می‌شود که برای ابزارهای گفتار در C به‌کار می‌رود.
۲. مدل زبان: مجموعه‌ای از آواها، قوانین گفتاری، آنالیز متن و مدل‌های نوایی را فراهم می‌آورد که اغلب در بازشناسی گفتار مورد استفاده قرار می‌گیرد. این قسمت از سیستم US نامیده می‌شود.
۳. مدل لغت‌نامه: یک مدل طرز تلفظ و ادای سخن بوده که شامل یک لغت‌نامه و قانون تبدیل حروف به صدا برای بیرون کشیدن الفاظ کلمات می‌باشد. لغت‌نامه‌ها به واحدهای ذخیره‌شده زبان وابسته هستند، این قسمت از سیستم CMU نامیده می‌شود. از نظر اندازه در نرم‌افزار Flite می‌توان گفت که لغت‌نامه‌ها و سیستم‌های قانون تبدیل حروف به صدا بیشترین حافظه نرم‌افزار را (در زبان انگلیسی) تشکیل می‌دهند.
۴. صدا: شامل فهرستی از مدل‌های نوایی یک گوینده خاص و تعریف صدای او می‌باشد. موارد ابتدایی در هر صدا توسط مدل زبان مربوط به آن، فراهم می‌شود. در این نرم‌افزار فایل صوتی با فرمت

پایه‌ریزی شده است [۹، ۲۰]. حداقل سه فرمنت برای تولید یک گفتار قابل مفهوم لازم است و اگر به پنج عدد برسد می‌توان گفتار باکیفیت بالا ایجاد کرد. هر فرمنت معمولاً با یک تشدیدکننده دوقطبی که قادر است هم فرکانس فرمنت و هم پهنای‌بندش را مشخص کند، مدل می‌شود.

## ۲-۳-۳ سنتز الحاقی

در دو روش قبل پارامترهای مشخصه گفتار در هر بازه زمانی توسط مجموعه‌ای از قواعد تولید می‌شدند، اما در این روش واحدهای گفتار ذخیره‌شده طبیعی برای تولید گفتار خروجی به‌صورت تکه‌تکه در کنار هم قرار می‌گیرند [۹، ۲۰]. یکی از مهم‌ترین جنبه‌ها در سنتز الحاقی انتخاب طول واحد صحیح است. دو مشکل در سنتز الحاقی نسبت به دیگر روش‌ها موجود است:

- بروز ناپیوستگی در نقاط اتصال که با استفاده از دوواجی‌ها و روش‌های خاصی برای هموارسازی سیگنال می‌توان آن را حل کرد.
- نیاز به حافظه بالا مخصوصاً زمانی که طول واحدها هم‌چون کلمه‌ها یا هجاها باشد.

## ۳- نحوه کار و تبدیل متن به گفتار در نرم‌افزار Flite

### ۳-۱ مقدمه

Flite که هسته اصلی آن از سیستم سنتزکننده گفتار Festival [۵] گرفته شده است، یک ماشین کوچک سنتزکننده گفتار زبان انگلیسی با زمان اجرای سریع است که برای سیستم‌های تعبیه‌یافته و سرورها مناسب می‌باشد [۴]. این ماشین که ابزارها، الگوها و مدارک موردنیاز برای ساخت یک صدای ساختگی (مصنوعی) جدید را فراهم می‌آورد، در ANCI C نوشته شده است و طوری طراحی شده که تقریباً به هر پایگاهی قابل‌انتقال بوده و دارای سخت‌افزار بسیار کوچکی است. درواقع Flite یک کتابخانه سنتز است که می‌تواند در داخل دیگر برنامه‌ها قرار گیرد.

### ۳-۲ سنتز قطعات گفتار در Flite

واحد گفتار در این سنتز کننده، الحاقی از نوع PSOLA TD<sup>۲۲</sup> است که درواقع الگوریتم اضافه کردن هم‌پوشانی گام‌های هم‌آهنگ در حوزه زمان می‌باشد.

در این روش از علائم گام<sup>۲۳</sup> به‌جای دوره تناوب گام استفاده می‌شود. علائم گام، گام‌هایی هستند که بزرگ‌ترین قله را در هر دوره تناوب دارند. برای تعیین علائم گام ابتدا شکل موج را توسط دو فیلتر بالاگذر و پایین‌گذر (یک فیلتر میان‌گذر) با توجه به جنسیت گوینده محدود می‌کنیم. در این نرم‌افزار به‌طور قراردادی حد پایین ۰/۰۰۵ ثانیه و حد بالا ۰/۰۱۲ ثانیه متعلق به یک گوینده مرد (محدوده فرکانسی ۸۰ تا ۲۰۰ هرتز) انتخاب شده است. با توجه به اینکه

- هر واج به یکی از گروه‌های واکه<sup>۲۵</sup> (+) یا هم‌خوان<sup>۲۶</sup> (-) تقسیم می‌شود. (VC<sup>TV</sup> + -)

**جدول ۱: حروف فارسی و نشانه‌ها و واج‌های معادل انگلیسی.**

حرف (نشانه)	حرف (نشانه) انگلیسی	واج	حرف (نشانه) فارسی	حرف (نشانه) انگلیسی	واج
آ	/aa/	/aa/	ص	/s/	/s/
ب	/b/	/b/	ض	/z/	/z/
پ	/p/	/p/	ط	/t/	/t/
ت	/t/	/t/	ظ	/z/	/z/
ث	/s/	/s/	ع	-	-
ج	/j/	/j/	غ	/gh/	/gh/
چ	/ch/	/ch/	ف	/f/	/f/
ح	/h/	/h/	ق	/gh/	/gh/
خ	/kh/	/kh/	ک	/k/	/k/
د	/d/	/d/	گ	/g/	/g/
ذ	/z/	/z/	ل	/l/	/l/
ر	/r/	/r/	م	/m/	/m/
ز	/z/	/z/	ن	/n/	/n/
ژ	/zh/	/zh/	و	/v/	/v/
س	/s/	/s/	ه	/hh/	/hh/
ش	/sh/	/sh/	ی	/y/	/y/

**جدول ۲: نشانه‌ها و واج‌های انگلیسی معادل ۶ واکه فارسی.**

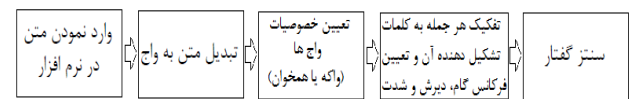
واج	نشانه انگلیسی	مصوت فارسی
/ae/, /aa/, /ax/	/a/	فتحه
/eh/	/e/	کسره
/ao/, /aa/, /ax/, /ow/	/o/	ضمه
/ae/, /aa/	/a:/	آ - ا
/ih/, /iy/, /ax/	/ee/, /i/	ی
/uw/, /uh/, /ow/, /ax/	/u:/	و

- در صورت واکه بودن:
  - به کدام یک از ۴ گروه زیر تعلق دارند؟
  - کوتاه (s)، بلند (l)، واکه دوگانه<sup>۲۸</sup> (d)، حرف صدادار میان کلمه بدون تشدید<sup>۲۹</sup> (a).
- در هنگام ادا کشیدگی صدا در چه حدی است: زیاد (۱)، متوسط (۲)، کم (۳)
- در هنگام ادا کشیدگی صدا در چه حدی است: زیاد (۱)، متوسط (۲)، کم (۳)
- در هنگام ادا کشیدگی صدا در چه حدی است: زیاد (۱)، متوسط (۲)، کم (۳)
- چگونه ادا می‌شوند. زبان با لب‌ها لمس می‌شود (+)، زبان گرد می‌شود (-)
- در صورت هم‌خوان بودن:
  - به کدام یک از گروه‌های زیر تعلق دارند؟

wav، با مشخصات مونو، ۱۶ بیت، ۸ کیلوهرتز pcm بوده و گوینده مرد می‌باشد.

**۴- مدل پیشنهادی برای تولید گفتار به زبان فارسی**

این تحقیق با استفاده از برنامه سنتزکننده گفتار جامع و کامل Flite انجام شده است. ابتدا به منبع برنامه دسترسی پیدا کرده و زیربرنامه‌های موردنیاز در قالب C در کنار هم متصل شده است. سپس با تغییرات انجام‌شده و اضافه نمودن اطلاعات موردنیاز بر روی این زیربرنامه‌ها، سنتزکننده گفتار در زبان فارسی تهیه شده است. در شکل ۲ روندنمای تبدیل متن به گفتار نرم‌افزار Flite نشان داده شده است که در قسمت مربوط به واج‌ها، اطلاعات فرکانس گام مربوط به زبان فارسی وارد شده است.



**شکل ۲: روندنمای تبدیل متن به گفتار در نرم‌افزار Flite**

بدین طریق که در بلوک وارد نمودن متن در نرم‌افزار، ابتدا در قسمت اصلی و فراخوان زیربرنامه، نام فایلی را که نرم‌افزار باید متن را از آن بخواند وارد می‌کنیم. همچنین نام فایلی (با فرمت Wave) را که می‌خواهیم متن خوانده‌شده در آن قرار گیرد، به برنامه اضافه شده است. در بلوک بعدی استخراج صورت واجی متن انجام می‌شود. برای مثال عبارت «دیروز برف بارید» به واج‌های زیر تبدیل می‌شود (pau به معنای وقفه می‌باشد):

"diruz barf barid"  
pau d ih r uw z b aa r f b aa r ih d pau

به‌طور کلی تبدیل حرف به صدا به دو دسته مبتنی بر فرهنگ لغت و مبتنی بر قانون طبقه‌بندی می‌شود. فرهنگ لغت این نرم‌افزار شامل ۱۰۰۰۰ کلمه می‌باشد. همچنین از قواعد تبدیل حروف به صدا نیز به‌صورت یک پشتیبان برای کلماتی که صریحاً در لغت‌نامه نیامده‌اند، استفاده می‌شود. نرم‌افزار علاوه بر یک فرهنگ‌نامه لغت، از فهرستی به نام افزایش نیز پشتیبانی می‌کند که برنامه‌های مشخص و کاربران آن‌ها را قادر می‌سازد تا ورودی‌های جدیدی را به لغت‌نامه اضافه کنند.

در جدول ۱ حروف فارسی و نشانه و واج‌های معادل انگلیسی نشان داده شده است. همچنین نشانه و واکه انگلیسی معادل ۶ واکه فارسی در جدول ۲ نشان داده شده است. در جدول ۳ تبدیل متن به واج ۴۰ جمله گرفته‌شده از پایگاه داده فارسی نشان داده شده است. بلوک بعدی مربوط به تعیین خصوصیات واج‌ها می‌باشد. در ادامه برنامه‌های مربوط به این بلوک، خصوصیات هر یک از واج‌ها در قالب نمودار درختی آورده شده است که به‌صورت زیر می‌باشد:

pau p aa eh z f ae s l a x b aa r g r ih z aa n ae s t pau	Paeez fasleh bargrizaan ast.	۲۳- پایین فصل برگ ریزان است.
pau d aa z sh eh k ax s t pau	Daas shekast.	۲۴- داس، شکست.
pau p aa r ch l ae b ax t ae g ch ax ae s t pau	Parch labezh tagcheh ast.	۲۵- پارچ لب طاقچه است.
pau s ae g aa l iy r aa g aa z g aa r eh f t pau	Sag ali ra gaaz gehref.	۲۶- سگ علی را گاز گرفت.
pau ih n z aa r f d aa r n aa d aa aa d pau	In zarf dar nadaarad?	۲۷- این ظرف، در ندارد؟
pau f ae g er b ae ey s ax y aa s ax m ih sh ax v ax d pau	Faghre baaish yase mishavad.	۲۸- فقر باعث یاس می‌شود.
pau jh ae n b ax k aa n iy ax m aa y eh k b aa ng k ae s t pau	Janbeh khaneh ma yak bank ast.	۲۹- جنب خانه ما یک بانک است.
pau m ao r g hh ae n uw z z eh n d ax ae s t pau	Morgh hanuz zende h ast.	۳۰- مرغ هنوز زنده است.
pau g aa v m aa d aa r ae m r aa sh aa k z ae d pau	Gaav maadaram ra shaakh zad.	۳۱- گاو، مادرم را شاخ زد.
pau ae s b f aa r aa r k aa d pau	Asbe faraar kard.	۳۲- اسب، فرار کرد.
pau r ae ng g ae ch z aa r d ae s t pau	Range gache zard ast.	۳۳- رنگ گچ زرد است.
pau v ae z ax k ae f sh aa ay ax t k ay l iy k eh r aa b ae s t pau	Vazeh kafshhaayat kheili kharaab ast.	۳۴- وضع کفش‌هایت خیلی خراب است.
pau hh ae n uw z y ae k d aa r iy m pau	Hanuz yakh daarim.	۳۵- هنوز یخ داریم.
pau b eh t ow ae s l a x n r ae b t n aa d aa aa d pau	Beh toe aslan rabt nadaarad.	۳۶- به تو اصلاً ربط ندارد.
pau ax jh ae b uw ax ae t r iy m iy aa ay ax d pau	Ajab booyeh aatari miaayad.	۳۷- عجب بوی عطری می‌آید.
pau g aa r k d ih r m ih p ax z ax d pau	Ghaaarch deer mipazad.	۳۸- قارچ دیر می‌پزد.
pau ae b g aa r m n ih s t pau	Ab garm nist?	۳۹- آب، گرم نیست؟
pau ih n r aa p ay aa n n aa d aa aa d pau	In raah paayaan nadaarad.	۴۰- این راه، پایان ندارد.

وقفه‌ها (s)، سایشی‌ها<sup>۳۳</sup> (f)، ترکیب وقفه‌ها و سایشی‌ها<sup>۳۴</sup> (a)، خیشومی‌ها<sup>۳۵</sup> (n)، شبه‌واکه مایع<sup>۳۶</sup> (l) (ctype<sup>37</sup> s f a n l . )  
 • در هنگام تلفظ کدام قسمت واج‌گاه مشغول می‌شود.  
 لب‌ها<sup>۳۸</sup> (l)، لب‌ها به‌صورت حباب‌دار جمع می‌شوند<sup>۳۹</sup> (a)، سخت‌کامی<sup>۴۰</sup> (p)، لب‌ها و دندان‌ها (b)، دندان‌ها (d)، نرم‌کامی<sup>۴۱</sup> (v). (cplace<sup>42</sup> l a p b d v . )  
 • واک‌دار یا بی‌واک هستند.

(cvox<sup>43</sup> + -)

جدول ۳: تبدیل متن به واج.

جمله فارسی	جمله معادل فارسی	تبدیل متن به واج
۱- دیروز برف بارید.	Dinuz barf barid.	Pau d ih r uw z b aa r f b aa r ih d pau
۲- من یک دستگاه چاپ ساختم.	Man yek datgaah chaap saakhtam.	Pau m ae n y eh k d ae s t g aa hh ax ch aa p s aa k t ax m pau
۳- پارچ توی کمد است.	Parch tooyeh komod ast.	Pau p aa r ch t uw ax k ow m aa d ae s t pau
۴- کتفم درد می‌کند.	Ket faam dard mikonad.	Pau k eh t f ae m d aa r d m ih k ax n ae d pau
۵- به چه علتی سگ فرار کرد؟	Beh che ellati sag faraar kard?	Pau b eh c hey l aa t iy s ae g faa r k aa d pau
۶- دیروز در کرج، قاسم را دیدم.	Deeruz dar karaj gaasem ra didam.	Pau d ih uw z d aa r k ae r ay g az s ax m r aa d ih d ax m pau
۷- گرگ به گله زد.	Gorge beh galleh zad.	Pau g ao r g b eh g ae l ax z ae d pau
۸- پایش از مچ شکست.	Payash az mohche shekast.	Pau p aa y ax sh ae z m ow ch sh eh k ax s t pau
۹- پاپ فردا می‌آید ایران.	Paap farada miaayad eeraan.	Pau p aa p f aa r d aa m iy aa ay ax d eh r aa n pau
۱۰- سوت گم شد.	Soot gome shod.	Pau s uh t g ow m sh ow d pau
۱۱- گاو توی مزرعه است.	Gaav tooyeh mazraeh ast.	Pau g aa v t uw ax m ael z r ey aa s t pau
۱۲- پیچ باز شد.	Pich baaz shod.	Pau p ih ch b aa z sh ow d pau
۱۳- جیب راه افتاد.	Jip rah ohftaad.	Pau jh ih p r aa hh aa l f t aa d pau
۱۴- آن اسب چه می خورد؟	An asb cheh mikhorad?	Pau ae n ae s b ch ehl m ih k hh ao r ax d pau
۱۵- پمپ، آب ندارد.	Pomp aab nadaarad.	Pau p aa m p ae b n aa d aa aa d pau
۱۶- دو رکعت، نماز خواند.	Doe rakat namaaz khaand.	Pua d ow r ae k ae t n aa m aa z k ax n d pau
۱۷- پول را قاب زد.	Pule ra gaap zadam.	Pau p uw l r aa g aa p z aa d ax m
۱۸- سرتیپ وقت ندارد.	Sartipe vagt nadaarad.	Pau s aa l t ax p v ae g t n aa d aa aal d pau
۱۹- دیروز توت چیدم.	Deeruz toot chidam.	Pau d ih r uw z t uw t chi h d ax m pau
۲۰- کتابت پاره شده است.	Kehtaabat paareh shodeh ast.	Pau k eh t ax b aa t p er ax sh owl eh ae s t pau
۲۱- من لوسش نکرده‌ام.	Man lusash nakradeham.	pau m ae n l uw s ax sh n ax k aa r d ax hh ax m pau
۲۲- حجاج به مکه رفتند.	Hojjaaj beh makkeh raftand.	pau hh oy jh aa jh b eh m ae k ax r ae f t ax n d pau

(phoneme	vc	vInlg	vheight	vfront	vrnd	ctype	cplace	cvox)
(#	-	1	-	-	-	.	.	-)
(a	+	1	۳	۱	-	.	.	-)
(e	+	1	۲	۱	-	.	.	-)
(i	+	1	۱	۱	-	.	.	-)
(o	+	1	۳	۳	-	.	.	-)
(u	+	1	۱	۳	+	.	.	-)
(b	-	.	-	-	+	s	l	+
(ch	-	.	-	-	+	a	a	-)
(d	-	.	-	-	+	s	a	+
(f	-	.	-	-	+	f	b	-)
(g	-	.	-	-	+	s	p	+
(j	-	.	-	-	+	l	a	+
(k	-	.	-	-	+	s	p	-)
(l	-	.	-	-	+	l	d	+
(ll	-	.	-	-	+	l	d	+
(m	-	.	-	-	+	n	l	+
(n	-	.	-	-	+	n	d	+
(ny	-	.	-	-	+	n	v	+
(p	-	.	-	-	+	s	l	-)
(r	-	.	-	-	+	l	p	+
(rr	-	.	-	-	+	l	p	+
(s	-	.	-	-	+	f	a	+
(t	-	.	-	-	+	s	t	+
(th	-	.	-	-	+	f	d	+
(x	-	.	-	-	+	a	a	-)

شکل ۳: نمودار درختی تابع Phonaset

(stddev) فرکانس‌های گام را برای هر واج واک‌دار جداگانه حساب می‌نماید.

- در مدل پیشنهادی برای شباهت بیشتر منحنی فرکانس گام به دست آمده از روش پیشنهادی با منحنی فرکانس گام پایگاه داده جملات فارسی فارس‌دات و نزدیک‌تر شدن صدای به دست آمده از مدل پیشنهادی به صدای انسان از سه معیار میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام هر واج استفاده می‌نماییم. یعنی علاوه بر وارد نمودن میانگین و انحراف معیار فرکانس‌های گام برای هر واج واک‌دار، بزرگ‌ترین فرکانس گام هر واج واک‌دار  $F_0$  را نیز به برنامه اضافه می‌نماییم.

- همان‌طور که گفتیم در رابطه ۱ مقدار ارزش ثابت در نظر گرفته شده است و با تغییر واج‌ها و تغییر میانگین و انحراف معیار تغییر نمی‌نماید؛ اما با تغییر رابطه ۱ به رابطه پیشنهادی ۲ مقدار ارزش را طوری تعیین نمودیم که حساس به تغییرات سه معیار میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام هر واج باشد و مقدار ارزش برای هر واج جداگانه تعیین می‌گردد.

در این قسمت منحنی گام مربوط به جمله «دیروز برف بارید» که از نرم‌افزار Flite به دست آورده‌ایم را با منحنی گامی که از پایگاه داده جملات فارسی فارس‌دات به دست آمده است، مقایسه می‌کنیم. همان‌طور که گفته شد، تغییرات فرکانس گام بر اساس میانگین و انحراف معیار فرکانس گام و  $F_0$  می‌باشد که از طریق محاسبه ارزش با توجه به رابطه پیشنهادی ۱ به نرم‌افزار وارد می‌شود. در جدول ۴ مقدار میانگین و انحراف معیار فرکانس گام هر واج واک‌دار مربوط به جمله «دیروز برف بارید» را با توجه به پایگاه داده جملات فارسی فارس‌دات، به دست آورده‌ایم. همچنین مقدار بزرگ‌ترین فرکانس گام مربوط به هر واج و مقدار ارزش نشان داده شده است.

در شکل ۴ سیگنال گفتار زمانی و منحنی فرکانس گام جمله «دیروز برف بارید» مربوط به نرم‌افزار Flite بدون اعمال روش پیشنهادی را می‌بینیم و در شکل ۵ سیگنال گفتار زمانی و منحنی فرکانس گام جمله «دیروز برف بارید» مربوط به نرم‌افزار Flite با اعمال روش پیشنهادی را می‌بینیم.

در شکل ۶ منحنی فرکانس گام پایگاه داده جملات فارسی فارس‌دات، نشان داده شده است. این نمودار را با دسترسی به اطلاعات فرکانس گام جملات فارسی از پایگاه داده فارس‌دات و از طریق نرم‌افزار چایلز که بر روی برنامه MATLAB ۵/۲ نصب شده است، ترسیم نموده‌ایم.

دلیل تفاوت پهنای پالس‌ها در منحنی فرکانس گام «دیروز برف بارید» مربوط به نرم‌افزار Flite با اعمال روش پیشنهادی با منحنی فرکانس گام پایگاه داده جملات فارسی مربوط به جمله «دیروز برف بارید» را می‌توان بدین صورت بیان نمود که در رسم منحنی فرکانس گام پایگاه داده جملات فارسی، تمام فرکانس‌های گام مربوط به هر واج واک‌دار در ایجاد منحنی دخیل هستند؛ اما در رسم منحنی فرکانس گام

شکل ۳ نمودار درختی تابع Phonetset را نشان می‌دهد. بلوک مربوط به عبارت‌بندی آوایی سخن را قابل‌فهم‌تر می‌کند. انسان بسته به اندازه شش‌ها می‌تواند مدت زمان محدودی قبل از نفس‌تازه‌کردن، صحبت کند. این موضوع یک حد بالا برای عبارت‌های آوایی تعیین می‌کند. با این حال، به ندرت عبارت‌هایمان را به طول بیشینه می‌سازیم. برای زبان انگلیسی (و احتمالاً بسیاری از زبان‌های دیگر) قواعد ساده بر پایه نقطه‌گذاری تعیین‌کننده مناسبی برای حدود عبارت‌های آوایی است. هم‌چنین، این حدود از طریق «ویرگول»، «علامت سؤال»، «علامت تعجب» تعیین می‌شوند. سپس در این بلوک فرکانس گام واج‌های صدادار (واک‌دار) را می‌آوریم.

تغییرات فرکانس گام بر اساس میانگین و انحراف معیار و بزرگ‌ترین فرکانس گام مربوط به هر واج واک‌دار می‌باشد که از طریق محاسبه ارزش مربوط به هر واج با توجه به رابطه ۱ در نرم‌افزار وارد می‌نماییم. برای این منظور به طریق زیر عمل می‌نماییم:

۱- میانگین (mean) و انحراف معیار (stddev) فرکانس‌های گام هر واج واک‌دار را به دست می‌آوریم.

۲- بزرگ‌ترین فرکانس گام هر واج واک‌دار را برابر  $F_0$  قرار می‌دهیم. در این برنامه مقادیر mm و ms بر اساس نرم‌افزار Flite ثابت می‌باشند.

۳- حال از طریق رابطه پیشنهادی ۲ مقدار ارزش را به دست می‌آوریم.

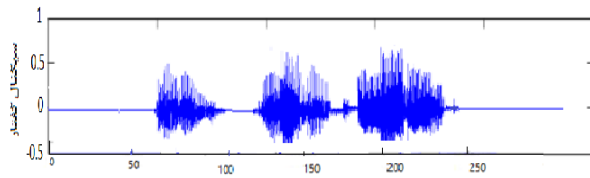
$$v = (((f_0 - m) / s) * ms) + mm \quad (2)$$

که در آن  $v = \text{value}$ ،  $m = \text{mean}$ ،  $s = \text{stddev}$ ،  $mm = \text{model\_mean} = 170$  و  $ms = \text{model\_stddev} = 34$  است.

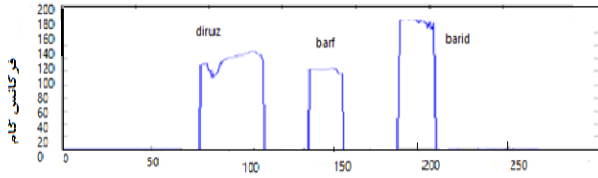
در رابطه پیشنهادی ۲ با استفاده از سه معیار میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام هر واج واک‌دار و گذاشتن در رابطه پیشنهادی، عددی به نام ارزش به دست می‌آوریم که این عدد تأثیر مستقیم در طبیعی‌تر و قابل‌فهم‌تر شدن صدای به دست آمده دارد. در واقع مقدار ارزش از طریق مقادیر میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام هر واج واک‌دار تعیین می‌گردد و مقدارش ثابت نمی‌باشد.

در برنامه نرم‌افزار Flite تنها می‌توانیم با تغییر میانگین و انحراف معیار که برای تمامی واج‌ها یکسان تعریف شده است، منحنی فرکانس گام تمام واج‌ها را به مقدار یکسان تغییر دهیم. همچنین مقادیری برای کل واج‌ها ثابت در نظر گرفته شده است که به نام ارزش (value)، مدل میانگین (model\_mean) و مدل انحراف معیار (model\_stddev) نام گذاری شده‌اند. در واقع این سه ثابت هیچ نقشی در فرکانس گام واج‌ها ندارند و برای تمامی واج‌ها یکسان تعریف شده‌اند. بنابراین ما نمی‌توانیم زیربومی هر واج را جداگانه تغییر دهیم. بنابراین صدای به دست آمده بسیار رباتیک می‌باشد و فقط با تغییر عدد مربوط به میانگین و انحراف معیار که برای کلیه واج‌ها می‌باشد، تغییرات یکسان به کار می‌رود؛ اما مدل پیشنهادی سه خصوصیت مهم دربر دارد:

- برنامه Flite را طوری تغییر دادیم تا بتوانیم زیربومی هر واج را جداگانه تغییر دهیم. بنابراین برنامه میانگین (mean) و انحراف معیار

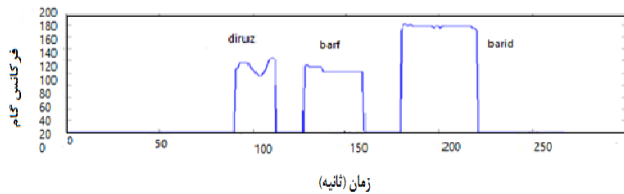


شکل ۴: جمله «دیروز برف بارید» بدون اعمال روش پیشنهادی (الف)



شکل ۵: جمله «دیروز برف بارید» با توجه به روش پیشنهادی (الف)

شکل ۵: جمله «دیروز برف بارید» با توجه به روش پیشنهادی (الف) سیگنال گفتار زمانی (ب) منحنی فرکانس گام، توسط نرم‌افزار Flite



شکل ۶: منحنی فرکانس گام پایگاه داده جملات فارسی مربوط به جمله «دیروز برف بارید»

### ۵- نتایج

برای ارزیابی سیستم‌های سنتز گفتار عمدتاً پارامترهای قابل فهم بودن، طبیعی بودن و خوش‌آیند بودن صدای سنتز شده، مورد ارزیابی قرار می‌گیرد. قابل فهم بودن، میزان وضوح و قابل درک بودن صدا را نشان می‌دهد. طبیعی بودن، مشخص می‌کند که صدای سنتز شده چقدر به صدای طبیعی انسان نزدیک است و خوش‌آیند بودن نیز مشخص می‌کند که صدای تولید شده تا چه حدی برای شنونده خوش‌آیند است. برای ارزیابی صدای سنتز شده، روش‌های مختلفی پیشنهاد شده است که آزمون MOS<sup>۴۴</sup> از جمله مهم‌ترین این روش‌ها به شمار می‌رود [۲۱]، [۲۲]. در آزمون MOS، چند جمله مختلف با مشخصات متفاوت از لحاظ گوناگونی واج‌ها، برای تعدادی شنونده پخش می‌شود. هر شنونده یک ارزیابی از پارامترهای قابل فهم بودن، طبیعی بودن و خوش‌آیند بودن ارائه می‌نماید. این ارزیابی به این صورت است که به ازای هر کدام از پارامترهای فوق یک نمره به هر جمله داده می‌شود. نمرات مورد استفاده و مفهوم آن‌ها، در جدول ۵ مشاهده می‌شود.

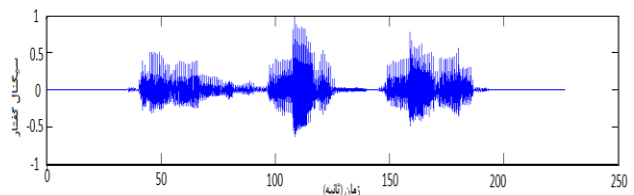
سپس به ازای هر کدام از پارامترهای فوق، از امتیازات میانگین‌گیری می‌کنیم و امتیاز هر پارامتر را مشخص می‌نماییم. هم‌چنین در این پروژه آزمون MOS را بر روی ۴۰ جمله از جملات موجود در پایگاه داده فارسی انجام دادیم که ۲۰ جمله اول مربوط به آموزش نرم‌افزار می‌باشند و ۲۰ جمله دوم برای آزمون نرم‌افزار به کار رفته است.

مربوط به نرم‌افزار Flite با اعمال روش پیشنهادی تنها از میانگین و انحراف معیار فرکانس‌های گام هر واج واک‌دار و بزرگ‌ترین فرکانس گام هر واج واک‌دار استفاده می‌نماییم.

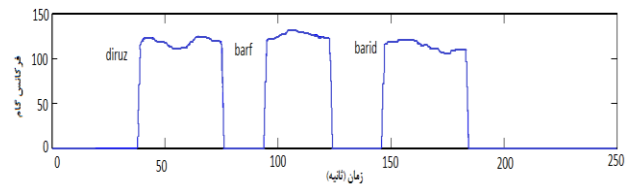
جدول ۴: تعیین میانگین، انحراف معیار، بزرگ‌ترین فرکانس گام و به‌دست آوردن مقدار ارزش مربوط به جمله «دیروز برف بارید».

	d	ih	r	uw	z	b	aa
واک‌دار	بی‌واک	واک‌دار	واک‌دار	واک‌دار	واک‌دار	بی‌واک	واک‌دار
نوع واج	---	۷۷	۱۱۱	۱۳۰	۲۱	---	۶۲
میانگین	---	۵۳	۸	۳	۴۷	---	۵۴
انحراف معیار	---	۱۹۶	۲۲۵	۲۱۵	۲۴۹	---	۲۰۱
ارزش	---	۱۱۸	۱۲۴	۱۳۴	۱۳۰	---	۱۱۲
بزرگ‌ترین فرکانس گام	---	---	---	---	---	---	---

	r	f	b	aa	r	ih	d
واک‌دار	واک‌دار	بی‌واک	واک‌دار	واک‌دار	واک‌دار	واک‌دار	بی‌واک
نوع واج	۱۰۵	---	۵۸	۱۷۴	۱۷	۱۰۳	---
میانگین	۳	---	۸۰	۴	۵۱	۷۴	---
انحراف معیار	۲۳۸	---	۲۲۰	۲۱۲	۲۷۳	۲۰۴	---
ارزش	۱۱۱	---	۱۸۰	۲۲۰	۱۷۲	۱۷۷	---
بزرگ‌ترین فرکانس گام	---	---	---	---	---	---	---



شکل ۴: جمله «دیروز برف بارید» بدون اعمال روش پیشنهادی (الف)



شکل ۴: جمله «دیروز برف بارید» بدون اعمال روش پیشنهادی (الف) سیگنال گفتار زمانی (ب) منحنی فرکانس گام، توسط نرم‌افزار Flite

شکل ۴: جمله «دیروز برف بارید» بدون اعمال روش پیشنهادی (الف) سیگنال گفتار زمانی (ب) منحنی فرکانس گام، توسط نرم‌افزار Flite

همچنین روش پیشنهادی بر روی ۴۰ جمله متفاوت (در جدول ۳ ذکر شده است)، انجام گرفته و از مقایسه منحنی فرکانس گام به‌دست آمده از روش پیشنهادی با منحنی فرکانس گام پایگاه داده جملات فارسی فارس دات، نتیجه مطلوب به‌دست آمده است.



جدول ۵: نمرات و مفهوم آن‌ها در تست MOS.

مفهوم	نمره
بد	۱
ضعیف	۲
نسبتاً خوب	۳
خوب	۴
فوق العاده خوب	۵

برای انجام آزمون MOS، از ۴۰ نفر شنونده استفاده نموده‌ایم. جدول ۶ نتایج حاصل از آزمون MOS، بر روی جملات آموزش با اعمال روش پیشنهادی و بدون اعمال روش پیشنهادی را نشان می‌دهد. همچنین جدول ۷ نتایج حاصل از آزمون MOS، بر روی جملات حاصل از آزمون نرم‌افزار با اعمال روش پیشنهادی و بدون اعمال روش پیشنهادی را نشان می‌دهد.

جدول ۶: مقایسه نتایج حاصل از آزمون MOS، بر روی جملات آموزش با توجه به اعمال روش پیشنهادی و بدون اعمال روش پیشنهادی بر

روی نرم‌افزار Flite.

جمله فارسی	قابل فهم بودن		طبیعی بودن		خوشایند بودن	
	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی
۱- دیروز برف بارید.	۳/۹	۳/۱	۳/۸	۳/۲	۴/۸	۳/۶
۲- من یک دستگاه چاپ ساخته‌ام.	۳/۸	۳/۸	۳/۹	۲/۴	۳/۸	۳/۵
۳- پارچ توی کمد است.	۴/۶	۳/۵	۴	۳/۶	۴/۳	۳/۲
۴- کتفم درد می‌کند.	۴/۵	۳/۲	۴/۱	۳/۵	۴/۱	۳/۱
۵- به چه علتی سگ فرار کرد؟	۴/۳	۳/۶	۳/۷	۳/۲	۳/۸	۳/۶
۶- دیروز در کرج، قاسم را دیدم.	۳/۹	۳/۸	۴/۵	۳/۴	۴/۳	۳/۷
۷- گرگ به گله زد.	۳/۷	۳/۳	۴/۴	۲/۳	۴/۵	۳/۵
۸- پایش از مچ شکست.	۴/۳	۳/۷	۳/۸	۳/۲	۴/۷	۳/۶
۹- پاپ فردا می‌آید ایران.	۴/۷	۳/۲	۴/۴	۳/۶	۴/۸	۳/۵
۱۰- سوت گم شد.	۴/۹	۳/۳	۴/۲	۳/۷	۴/۵	۳/۲
۱۱- گاو توی مزرعه است.	۴/۴	۳/۸	۳/۸	۳/۱	۴/۹	۳/۵
۱۲- بیج باز شد.	۴/۸	۳/۷	۴/۲	۳/۲	۴/۷	۳/۲
۱۳- جیب راه افتاد.	۳/۷	۳/۲	۴/۶	۳/۷	۴/۹	۳/۵
۱۴- آن اسب چه می‌خورد؟	۴/۵	۳/۱	۴/۲	۳/۵	۴/۷	۳/۳
۱۵- پمپ، آب ندارد.	۴/۱	۳/۲	۳/۸	۳/۲	۴/۹	۳/۵
۱۶- دو رکعت، نماز خواند.	۴/۵	۳/۴	۴/۵	۳/۷	۴/۷	۳/۵
۱۷- پول را قاب زد.	۴/۷	۳/۲	۴/۵	۳/۴	۴/۸	۳/۵
۱۸- سرتیپ وقت ندارد.	۴/۴	۳/۳	۴/۳	۳/۸	۴/۶	۳/۲
۱۹- دیروز توت چیدم.	۴/۶	۳/۱	۳/۹	۳/۳	۴/۷	۳/۵
۲۰- کتابت پاره شد.	۴/۹	۳/۲	۴/۷	۳/۲	۴/۵	۳
میانگین	۴/۴	۳/۳	۴/۲	۳/۴	۴/۶	۳/۴

جدول ۷: نتایج حاصل از آزمون MOS، بر روی جملات آزمون با توجه به اعمال روش پیشنهادی و بدون اعمال روش پیشنهادی بر روی

نرم‌افزار Flite.

جمله فارسی	قابل فهم بودن		طبیعی بودن		خوشایند بودن	
	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی	بدون اعمال روش پیشنهادی
۲۱- من لوسش نکرده‌ام.	۳/۹	۳/۳	۴	۳/۶	۴/۱	۳/۷
۲۲- حجاج به مکه رفتند.	۳/۷	۳/۴	۳/۹	۳/۵	۳/۹	۳/۴
۲۳- پاییز فصل برگ ریزان است.	۴/۱	۳/۲	۳/۹	۳/۴	۴/۸	۳
۲۴- داس، شکست.	۴/۳	۳/۲	۴/۲	۳/۵	۴/۳	۳/۲
۲۵- پارچ لب طاقچه است.	۴/۶	۳/۶	۴/۵	۳/۳	۴/۴	۳/۳
۲۶- سگ علی را گاز گرفت.	۴/۵	۳/۳	۳/۶	۳/۲	۴/۶	۳/۴
۲۷- این ظرف، در ندارد؟	۴/۲	۳/۵	۴/۲	۳/۵	۳/۸	۳/۵
۲۸- فقر باعث یاس می‌شود.	۴/۴	۳/۲	۴/۴	۳/۴	۳/۷	۳/۱
۲۹- جنب خانه ما یک بانک است.	۳/۹	۲/۸	۴/۶	۳/۲	۳/۹	۳
۳۰- مرغ هنوز زنده است.	۴/۴	۳/۲	۳/۷	۳/۳	۴/۵	۳/۵
۳۱- گاو، مادرم را شاخ زد.	۳/۹	۳/۲	۳/۸	۳/۵	۳/۸	۳/۴
۳۲- اسب، فرار کرد.	۳/۹	۳/۳	۴/۵	۳/۶	۴/۹	۳/۱
۳۳- رنگ گچ زرد است.	۳/۷	۳/۴	۴/۳	۳/۶	۴/۲	۳/۴
۳۴- وضع کفش هایت خیلی خراب است.	۴/۱	۳/۲	۳/۶	۳/۱	۴/۶	۳/۶
۳۵- هنوز بیخ داریم.	۴/۷	۳/۲	۳/۸	۳/۴	۳/۹	۳/۱
۳۶- به تو اصلاً ربط ندارد.	۳/۸	۳/۴	۴/۶	۳	۴/۴	۳/۷
۳۷- عجب بوی عطری می‌آید.	۴/۳	۳/۳	۳/۹	۳/۳	۴/۵	۳/۲
۳۸- قارچ دیر می‌پزد.	۴/۴	۳/۴	۴	۳/۵	۳/۸	۳/۴
۳۹- آب، گرم نیست؟	۳/۹	۳/۴	۳/۹	۳/۲	۳/۷	۳/۲
۴۰- این راه، پایان ندارد.	۴/۵	۳/۶	۴/۶	۳/۲	۴/۸	۳/۵
میانگین	۴/۲	۳/۳۹	۴/۱	۳/۳۶۵	۴/۳	۳/۳۳۵

**جدول ۸: مقایسه نتایج حاصل از آزمون MOS، بر روی سنتز جملات آزمون روش پیشنهادی با روش‌های مشابه پیشین.**

قابل فهم بودن	طبیعی بودن	خوشایند بودن	
۴/۴	۴/۲	۴/۸	جملات سنتز شده با استفاده از روش پیشنهادی اعمال شده بر Flite روی نرم افزار
۳/۳	۳/۴	۳/۴	جملات سنتز شده بدون اعمال روش پیشنهادی روی نرم افزار
۳/۸	۳/۹	۲/۵	جملات سنتز شده با استفاده از مدل‌های تراپون برای طول واج و گام [۶]
۴/۲	۴/۴	۴/۱	جملات سنتز شده با استفاده از درخت تصمیم گیری برای طول زمانی واج و گام [۶]
۴/۶	۴/۳	۴/۵	جملات سنتز شده با استفاده از درخت تصمیم گیری برای طول زمانی واج‌ها و روش اتوماتیک برای تولید گام [۶]
۴/۱۸	۳/۴۲	۳	سنتز کننده پیاده سازی شده به روش هارمونیک نوبزی [۷]

در مدل پیشنهادی ابتدا با توجه به پایگاه داده جملات فارسی فارسی‌دات، میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام هر واج واک‌دار را به دست می‌آوریم، سپس تغییرات هر واج عبارت موردنظر را از طریق محاسبه مقدار ارزش آن با توجه به رابطه ارائه شده در بخش روش پیشنهادی، در نرم‌افزار وارد می‌کنیم تا مشکل اول حل شود و منحنی گام هر واج با توجه به مقدار میانگین، انحراف معیار و بزرگ‌ترین فرکانس گام خودش تغییر کند، همچنین در این روش فقط برای واج‌های واک‌دار فرکانس گام در نظر گرفتیم و برای واج‌های بی‌واک، میانگین و انحراف معیاری در نظر گرفته نشد تا واج‌های واک‌دار از بی‌واک تشخیص داده شوند و مشکل دوم نیز حل شود. نتایج حاصل از آزمون‌های شنیداری، برای میزان قابل فهم بودن، طبیعی بودن و خوش آیند بودن برای جملات آموزش به ترتیب ۴/۴، ۴/۲ و ۴/۶ می‌باشد. همچنین، نتایج برای جملات مجموعه آزمون، به ترتیب برابر با ۴/۲، ۴/۱ و ۴/۳ می‌باشد.

## مراجع

- [1] A. B. Black and K. A. Lenzo, *Building synthetic voices*, For FestVox 2.1 Edition, 2007.
- [2] محمدمهدی همایون پور، محمد ایزدی، «تبدیل حرف به صدا در سیستم‌های تبدیل متن به گفتار فارسی با استفاده از درخت‌های تصمیم‌گیری CART»، *دوازدهمین کنفرانس سالانه انجمن کامپیوتر ایران*، تهران، ۱۳۸۵.
- [3] R. Kurzweil, *The singularity is near*, Penguin Books, ISBN 9-303788-14-0, 2005.
- [4] A.B. Black and K.A. Lenzo, *Flite: a small, fast speech synthesis engine*, System documentation Edition 1.3, for Flite version 1.3, 2005.
- [5] A. B. Black, P. Taylor and R. Caley, *The Festival Speech Synthesis System*, [Online], Available: <http://www.cstr.ed.ac.uk/projects/festival.-html>, 1998.
- [6] محمدمهدی همایون پور، مجید نم‌نبات، «تبدیل حرف به صدا در زبان فارسی به کمک شبکه‌های عصبی پرسپترون چندلایه‌ای»، *فصلنامه مهندسی برق و مهندسی کامپیوتر ایران*، شماره ۳، صفحات ۱۴۷-۱۵۴، پائیز ۱۳۸۶.
- [7] محمدمهدی همایون پور، سیدمصطفی موسوی، «تولید پارامترهای سنتز گفتار فارسی با استفاده از مدل‌های مخفی مارکوف و درخت تصمیم‌گیری»، *نشریه علمی - پژوهشی انجمن کامپیوتر ایران*، شماره ۱ و ۳ (الف)، صفحات ۱۹-۳۰، بهار و پائیز ۱۳۸۳.
- [8] Y. Sagisaka, "Speech synthesis from text," *IEEE Commun. Mag.*, pp. 35-41, 1990.
- [9] منصور شیخان نصیرزاده. مجید و دفتریان. علی، «طراحی و پیاده‌سازی سیستم تبدیل متن به گفتار طبیعی برای زبان فارسی»، *مجله علمی- پژوهشی دانشکده مهندسی دانشگاه فردوسی مشهد*، شماره ۲، صفحات ۳۱-۴۸، ۱۳۸۴.
- [10] محمدمهدی همایون پور، آرمین سلیمی بدر، «تعیین مرز و نوع عبارات نحوی در متون فارسی»، *فصلنامه علمی - پژوهشی پردازش علائم و داده‌ها*، شماره ۲، صفحات ۶۹-۸۶، ۱۳۹۲.
- [11] N. Thorensen, "Sentence intonation in textual context-supplementary data," *J. Acoust. Soc. Am.*, vol. 80, no. 4, pp. 1041-1047, 1986.

نتایج به دست آمده توسط آزمون MOS، برای پارامترهای قابل فهم بودن، طبیعی بودن و خوش آیند بودن برای جملات آموزش به ترتیب ۴/۴، ۴/۲ و ۴/۶ می‌باشند. همچنین برای جملات جهت آزمون نرم‌افزار، پارامترهای قابل فهم بودن، طبیعی بودن و خوش آیند بودن به ترتیب برابر ۴/۲، ۴/۱ و ۴/۳ می‌باشند. مقایسه این نتایج با حاصل روش‌های مشابه پیشین در جدول ۸ آورده شده است.

## ۶- نتیجه‌گیری

در این پژوهش، سنتز گفتار به زبان فارسی با استفاده از نرم‌افزار Flite مورد بررسی قرار گرفت. تغییرات منحنی گام در نرم‌افزار Flite با استفاده از تغییر میانگین و انحراف معیار فرکانس گام می‌باشد که دارای دو مشکل است:

برای تمام واج‌ها، یک میانگین و انحراف معیار واحد تعریف شده است یعنی با تغییر میانگین و انحراف معیار تعریف شده، منحنی گام تمام واج‌ها به یک اندازه تغییر می‌کند. همچنین در این نرم‌افزار برای صداهای بی‌واک نیز فرکانس گام در نظر گرفته شده است. بدین طریق که به وسیله درون‌یابی صداهای واک‌دار مجاور، به صدای بی‌واک (غیر از سکوت) فرکانس گام اختصاص داده و نتیجه آن هموارسازی می‌شود؛ بنابراین تشخیص صداهای واک‌دار از بی‌واک وجود ندارد.

- [18] P. H. Low and S. Vaseghi, "Application of microprosody models in TTS synthesis," *In Proc. ICSLP*, pp. 2413-2416, USA, 2002.
- [19] Y. Hifny. and M. Rashwan, "Duration modeling for Arabic TTS synthesis," *In Proc. ICSLP*, pp. 1773-1776, 2002.
- [20] A. Breen, "Speech synthesis models: a review," *Elect. Commun. Engng. J.*, pp. 19-31, 1992.
- [۲۱] سیدسعید آیت، طراحی و پیاده‌سازی سیستم تولید گفتار فارسی با تأکید بر بهبود هرچه بیشتر کیفیت گفتار تولیدشده، پایان‌نامه کارشناسی ارشد، دانشگاه صنعتی امیرکبیر، ۱۳۷۹.
- [22] S. Lemmetty, *Review of Speech Synthesis Technology*, Master Thesis, Helsinki University of Technology, 1999.
- [12] Y. Sagisaka, "On the prediction of global  $F_0$  shape for Japanese TTS," *In Proc. ICASSP*, USA, pp. 325-328, 1990.
- [13] J. Buhmann, et al. "Intonation modeling for the synthesis of structured documents," *In Proc. ICSLP*, USA, pp. 2089-2092, 2002.
- [14] M. Riedi, "A neural-network-based model of segmental duration for speech synthesis," *In Proc. Eurospeech*, Spain, pp. 599-602, 1995.
- [15] Z. Yiqing, "Syllable duration and its functions in standard Chinese discourse," *In Proc. ICSLP*, p. 1097, China, 2000.
- [16] C.L. Smith, "Modeling durational variability in reading aloud a connected text," *In Proc. ICSLP*, pp. 1769-1772, USA, 2002.
- [17] Y. Sagisaka, and Sato H. "Accentuation rules in Japanese TTS conversion," *Rev. Elect. Commun. Lab.*, vol. 32, no. 2, pp. 188-199, 1984.

## زیر نویس‌ها

<sup>1</sup> FarsDat: (Farsi Speech Database) [www.dadegan.ir](http://www.dadegan.ir)

<sup>2</sup> Stddev (Standard deviation)

<sup>3</sup> Phoneme

<sup>4</sup> Value

<sup>5</sup> Source- filter

<sup>6</sup> Concatenative

<sup>7</sup> LPC (Linear predictive coding)

<sup>^</sup> TTS (Text to speech)

<sup>9</sup> NLP (Natural Language Processor)

<sup>10</sup> Morphological

<sup>11</sup> POS (Part-Of-Speech tagging)

<sup>12</sup> Prosody

<sup>13</sup> Voiced

<sup>14</sup> Unvoiced

<sup>15</sup> Duration

<sup>16</sup> Intensity

<sup>17</sup> Pause

<sup>18</sup> Articulatory Synthesis

<sup>19</sup> Formant Synthesis

<sup>20</sup> Concatenative Synthesis

<sup>21</sup> Vocal tract

<sup>22</sup> Time-Domain Pitch Synchronous OverLap Add

<sup>23</sup> Pitch marks

<sup>24</sup> Residual

<sup>25</sup> Vowel

<sup>26</sup> Consonant

<sup>27</sup> Vowel or Consonant

<sup>28</sup> Diphthong

<sup>29</sup> Schwa

<sup>30</sup> Vowel length

<sup>31</sup> Vowel frontless

<sup>32</sup> Lip rounding

<sup>33</sup> Fricative

<sup>34</sup> Affricative

<sup>35</sup> Nasal

<sup>36</sup> Semi\_vowel liquid

<sup>37</sup> Vonsonant type

<sup>38</sup> Labial

<sup>39</sup> Alveolar

<sup>40</sup> Palatal

<sup>41</sup> Velar

<sup>42</sup> Place of articulation

<sup>43</sup> Consonant or Voicing

<sup>44</sup> Mean opinion score