

# تخمین مکان نواحی کدکننده پروتئین در توالی عددی DNA با استفاده پنجره با طول متغیر بر مبنای منحنی سه بعدی Z

حمیدرضا صابرکاری\*، موسی شمسی و محمد حسین صدیقی

دانشکده مهندسی برق، دانشگاه صنعتی سهند تبریز، تبریز، ایران



## چکیده

تخمین دقیق نواحی کدکننده پروتئین در ژن‌ها با استفاده از ابزارهای پردازش سیگنال، در سال‌های اخیر به چالشی در بیوانفورماتیک تبدیل شده است. بسیاری از روش‌های پردازش سیگنال‌های ژنومیک بر مبنای خاصیت تناوب-۳ بازهای موجود در رشته‌های DNA متمرکز بوده و سپس تحلیل‌های طیفی به منظور یافتن موقعیت مؤلفه‌های متناوب بر روی توالی‌های عددی DNA اعمال می‌شود. در این مقاله با استفاده از پنجره با طول متغیر و بر مبنای منحنی Z، الگوریتمی به منظور تعیین نواحی کدکننده پروتئین ارائه می‌کنیم. منحنی Z، یک منحنی سه بعدی منحصربه‌فرد برای نمایش توالی DNA می‌باشد که توصیف کاملی از رفتار بیولوژیکی توالی DNA را به دست می‌دهد. الگوریتم پیشنهادی به دلیل استفاده از پنجره گوسی با طول قابل تنظیم، از وضوح و دقت بسیار بالایی در تخمین نواحی ژنی برخوردار بوده و نواحی غیرپروتئینی در آن به طور کامل حذف می‌شود. همچنین به منظور استخراج مؤلفه تناوب-۳ از یک فیلتر میان‌گذر باند-محدود با فرکانس مرکزی  $\frac{2\pi}{3}$  استفاده می‌کنیم. الگوریتم پیشنهادی ابتدا بر روی توالی F56F11.4 در C.elegans اعمال و نتایج آن با سایر روش‌های موجود مقایسه شده و سپس آن را به ترتیب بر روی ژن‌های موجود در دو پایگاه داده HMR195 و BG570 اعمال می‌کنیم.

واژگان کلیدی: نواحی کدکننده پروتئین، تناوب-۳، DNA، پنجره با طول متغیر، فیلتر باند-محدود.

## Estimation of protein-coding regions in numerical DNA sequences using variable length window method based on 3D Z-curve

Hamidreza Saberkeri\*, Mousa Shamsi & Mohammad Hossein Sedaaghi

Department of Electrical Engineering, Sahand University of Technology, Tabriz, Iran

### Abstract

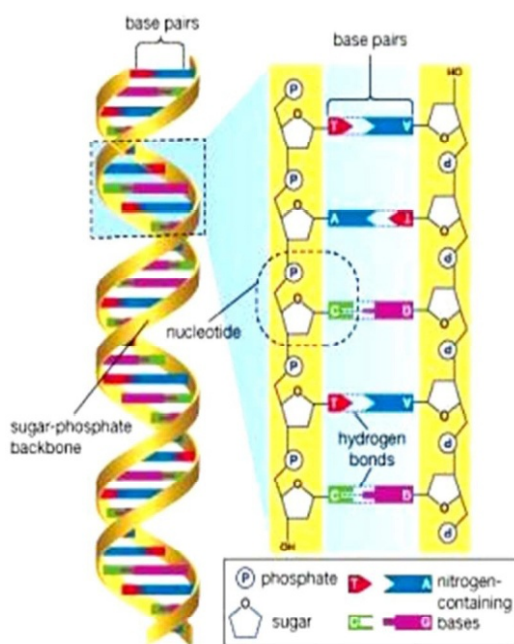
In recent years, estimation of protein-coding regions in numerical deoxyribonucleic acid (DNA) sequences using signal processing tools has been a challenging issue in bioinformatics, owing to their 3-base periodicity. Several digital signal processing (DSP) tools have been applied in order to identify the task and concentrated on assigning numerical values to the symbolic DNA sequence, then applying spectral analysis tools such as the discrete Fourier transform (DFT) to locate the periodicity components. Despite of many advantages of Fourier transform in detection of exotic regions, this approach has some restrictions, such as high computational complexity and disability in locating the small length coding regions. In this paper, we improve the performance of the conventional DFT in estimating the protein coding regions utilizing a

Gaussian window with variable length. First, the DNA strands are converted to numerical signals via the 3-D Z-curve method. Z curve is a robust, independent, less redundant approach, and has clear biological interpretation which can be regarded as a useful visualization technique for DNA analysis of any length. In the second stage, non-coding regions besides the background noise components are completely suppressed using the Gaussian variable length window. Also, we use a narrow-band band-pass filter in order to extract the period-3 components with  $\frac{2\pi}{3}$  central frequency. Performance of the proposed algorithm is tested on

F56F11.4 from C.elegans chromosome III, also two eukaryotic datasets, HMR195 and BG570, is compared with other state-of-the-art methods based on the nucleotide evaluation metrics such as sensitivity, specificity, approximation correlation, and precision. Results revealed that, the area under the receiver operating characteristic (ROC) curve is improved from 4% to 40%, in HMR195 and BG570 datasets compared to other methods. Furthermore, the proposed algorithm reduces the number of incorrect nucleotides which are estimated as coding regions.

**Keywords:** Protein coding regions, Period-3, Digital signal processing, DNA, Variable length window, Band-limited filter.

چون جهت‌های عکس هم دارند، موازی معکوس هستند. علاوه بر این، رشته‌های DNA جهت‌دار هستند (از انتهای دارای کربن 5' به طرف انتهای دارای کربن 3')، زیرا اتصال به نوکلئوتیدهای جدید همواره از جهت کربن 3' انجام یافته و در نتیجه تشکیل توالی جدید DNA در جهت 5' به 3' انجام می‌گیرد. همچنین به هنگام رمزدهی ساخت پروتئین، رشته DNA در این جهت پویش می‌شود.



(شکل-1): ساختار ملکول DNA [2]  
(Figure-1): DNA molecule structure [2]

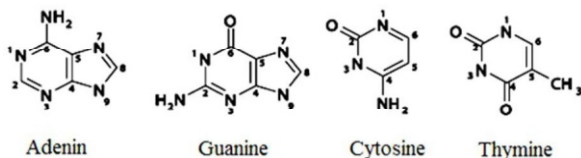
## ۱- مقدمه

هدف اصلی علم ژنتیک، درک ماهیت اطلاعات نهفته در ژن‌ها و نقش آن‌ها در تعیین عملکردهای خاصی است که توسط ژن‌ها بیان می‌شود. یکی از گام‌های اصلی برای رسیدن به این هدف، تعیین موقعیت ژن‌ها و به‌طور دقیق‌تر، تعیین نواحی کدکننده پروتئین در رشته‌های DNA<sup>۱</sup> است. یک رشته DNA، ملکول طولی از خانواده بیوپلیمرهاست که اطلاعات ژنتیکی را حمل می‌کند و دارای وظایف بیولوژیکی مهمی از جمله ذخیره و انتقال اطلاعات ژنتیک است. این ملکول، از دو رشته پلیمر خطی تشکیل شده و متشکل از واحدهای منومری به‌نام نوکلئوتید است [1]. شکل ۱ ساختار ملکول DNA را نشان می‌دهد. همان‌طور که مشاهده می‌شود، هر ملکول، حاوی نوکلئوتیدهای مختلف بوده و هر نوکلئوتید از سه بخش تشکیل شده است: یک قند پنج کربنه (قند دئوکسی ریبوز)، یک تا سه گروه فسفات ( $PO_4^-$ ) و نیز یک باز آلی نیتروژن‌دار. در شکل (۲)، ساختار یک نوکلئوتید در توالی DNA نشان داده شده است [2].

بازهای موجود در نوکلئوتیدها چهار نوع مختلف دارند؛ آدنین (A) و گوانین (G)، از دسته پورین‌ها بوده و دو حلقه‌ای هستند و تیمین (T)، سیتوزین (C) از دسته پیریمیدین‌ها بوده و تک‌حلقه‌ای هستند. شکل (۳) ساختار بازهای موجود در ملکول DNA را نشان می‌دهد [2]. در DNA بازهای مکمل (آدنین با تیمین و سیتوزین با گوانین) دوبده‌و با پیوند هیدروژنی به‌هم متصل شده و تشکیل جفت باز<sup>۲</sup> می‌دهند [2]. بنابراین، دو رشته DNA مکمل هم بوده و

<sup>1</sup> Deoxyribonucleic Acid

<sup>2</sup> Base-Pair



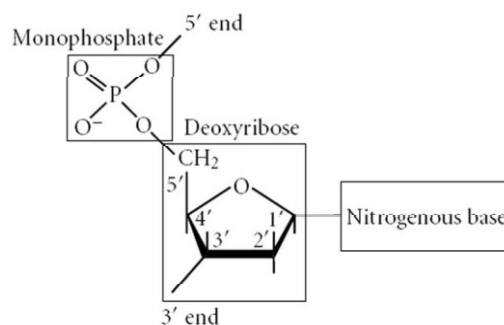
(شکل-۳): ساختار بازهای موجود در ملکول DNA [2]  
(Figure-3): Structure of bases in DNA molecule [2]

- اکسون داخلی<sup>۶</sup> که از یک ناحیه پذیرنده پروتئین شروع و تا نخستین ناحیه دهنده پروتئین بعدی ادامه می‌یابد.

- اکسون انتهایی<sup>۷</sup> که از آخرین ناحیه پذیرنده پروتئین شروع شده و تا کدون انتهایی (یکی از کدون‌های TAA، TAG و TGA) ادامه می‌یابد.

نواحی کد کننده پروتئین دارای مشخصه تناوب-۳<sup>۸</sup> هستند که در دیگر قسمت‌های ملکول DNA مشاهده نمی‌شود [5]. این پدیده می‌تواند به علت استفاده غیرهمگن از کدون‌ها باشد. این بدین معنی است که با وجود اینکه چند کدون ممکن است یک اسید آمینه خاص را رمزدهی کنند، همه آن‌ها با احتمال یکسانی در موجودات زنده ظاهر نمی‌شوند. به عنوان نوکلئوتید G در موقعیت‌های خاصی در کدون‌های نواحی اکسون جا می‌گیرد [7]-[6]. مشخصه تناوب-۳ می‌تواند به عنوان شناساگری برای تشخیص نواحی ژنی بکار رود.

حضور نوفه پس زمینه در توالی DNA یک مشکل اساسی در تعیین نواحی ژنی به شمار می‌آید [9]-[8]. همچنین به دلیل ماهیت پیچیده این نواحی، به طور معمول نیاز به یک ابزار قدرتمندی است که بتواند به طور مؤثر مشخصات نواحی کد کننده پروتئین را نمایش دهد. تاکنون روش‌های مختلفی برای رفع این مشکل پیشنهاد شده که در یک تقسیم‌بندی کلی می‌توان آن‌ها را به دو دسته تقسیم بندی کرد: الگوریتم‌های وابسته به مدل<sup>۹</sup> یا روش‌های نظارتی<sup>۱۰</sup> و الگوریتم‌های مستقل از مدل یا روش‌های مبتنی بر فیلتر<sup>۱۱</sup>. روش‌های وابسته به مدل مثل مدل مخفی مارکوف<sup>۱۲</sup> [10] و شبکه‌های عصبی<sup>۱۳</sup> [11]، که بر اساس برخی از اطلاعات اولیه جمع‌آوری شده از داده‌های موجود



(شکل-۲): ساختار یک نوکلئوتید در ملکول DNA  
(Figure-2): Structure of a nucleotide in DNA molecule

هر سه نوکلئوتید متوالی در توالی DNA، کدون<sup>۱</sup> نامیده می‌شود. هر کدون یک اسید آمینه خاص را فرا می‌خواند. پروتئین‌ها به عنوان واحدهای اساسی در بسیاری از فرآیندهای بیولوژیکی در داخل سلول، بافت و یا ارگانیسم‌های زنده، از آمینواسیدها تشکیل می‌شوند. سنتز پروتئین‌ها با استفاده از کدهای ژنتیکی صورت می‌گیرد؛ به گونه‌ای که ۶۴ کدون در توالی DNA به ۲۰ آمینواسید متناظر نگاشت می‌یابند؛ پس نگاشت از کدون‌ها به اسید آمینه‌ها یک نگاشت چند به یک است؛ یعنی ممکن است هر اسید آمینه توسط یک یا چند کدون فراخوانده شود. سه کدون نیز به عنوان کدون‌های توقف<sup>۲</sup> (کدون‌های TAA، TAG و TGA) وجود دارند که نشان دهنده پایان فرآیند پروتئین سازی است [3].

در یوکاریوت‌ها، DNA به دو ناحیه ژنی و بین ژنی تقسیم می‌شود. تنها ناحیه ژنی یا به اختصار ژن، اطلاعات را برای سنتز پروتئین‌ها حمل می‌کند. هر ژن نیز به نوبه خود متشکل از نواحی اکسون<sup>۳</sup> و انترون<sup>۴</sup> است که در شکل (۴) نشان داده شده است. بنابراین، اکسون‌ها کدهای لازم را برای تولید پروتئین حمل می‌کنند؛ از این رو به آن‌ها نواحی کد کننده پروتئین می‌گویند [4]. با مشاهده از انتهای 5' از توالی DNA به طرف انتهای 3' می‌توان دریافت به طور کلی سه ناحیه کد کننده پروتئین در توالی DNA وجود دارند (شکل ۴) که عبارتند از:

- اکسون ابتدایی<sup>۵</sup> که از کدون آغازی (کدون ATG) شروع و تا نخستین ناحیه دهنده پروتئین ادامه دارد.

<sup>6</sup> Donor Site

<sup>7</sup> Acceptor Site

<sup>8</sup> Period-3 Property

<sup>9</sup> Model-Dependent

<sup>10</sup> Supervised Methods

<sup>11</sup> Filter-based Methods

<sup>12</sup> Hidden Markov Model (HMM)

<sup>13</sup> Neural Networks

<sup>1</sup> Codon

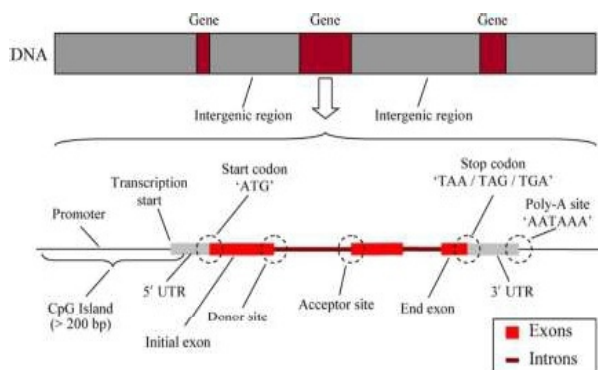
<sup>2</sup> Stop Codons

<sup>3</sup> Exon

<sup>4</sup> Intron

<sup>5</sup> Initial Exon

می‌شود. معیارهای ارزیابی سطح نوکلئوتیدی به‌منظور مقایسه الگوریتم پیشنهادی و سایر روش‌های موجود در بخش ۶ بیان می‌شود. نتایج پیاده‌سازی در محیط بسته نرم‌افزاری طراحی‌شده در بخش ۷ ارائه شده و درنهایت بخش ۸ شامل جمع‌بندی و نتیجه‌گیری از مقاله است.



(شکل-۴): نواحی اکسون و انترون در DNA یوکاریوتی [17]  
(Figure-4): Exon/Intron regions in Eukaryotic DNA [17]

## ۲- تبدیل رشته‌های DNA به سیگنال‌های دیجیتال

تبدیل رشته‌های DNA به سیگنال‌های دیجیتال، امکان اعمال تکنیک‌های پردازش سیگنال را به‌منظور تحلیل داده‌های ژنومیک و همچنین آشکارسازی ویژگی‌های کروموزوم‌ها نتیجه می‌دهد. ابزارهای تحلیل سیگنال‌های ژنومیک در حال حاضر توانایی آشکارسازی سطح وسیعی از ویژگی‌های یک دنباله DNA را در فواصل بیش از  $10^6 - 10^8$  جفت باز که شامل هر دو ناحیه پروتئینی و غیرپروتئینی است، در کلیه ژنوم‌ها یا کروموزوم‌ها به‌دست می‌دهند [19]-[20].

با اجرای هرچه بیشتر پروژه‌های ژنی، تعداد زیادی از توالی‌های ژنی در پایگاه‌های عمومی داده در دسترس است که نیاز ضروری به یافتن راه‌حل‌های جدید ریاضی به‌منظور تحلیل این توالی‌ها را نتیجه می‌دهد [21]. به‌طور کلی، دو روش پایه ریاضی در مطالعه تئوری‌های فیزیکی وجود دارد؛ یک روش، روش جبری و دیگری روش هندسی است. در بیش‌تر موارد این دو روش مکمل یکدیگر هستند. امروزه، در زمینه مطالعه ژن‌ها، روش جبری در تحلیل توالی ژنی به‌طور گسترده‌ای استفاده شده است، درحالی‌که روش هندسی برای تحلیل توالی ژنی مدت‌زمان زیادی است که صرف‌نظر شده است. منحنی  $Z$  یک منحنی فضایی سه‌بعدی است که به‌منظور نمایش توالی DNA به‌کار می‌رود. بر اساس این منحنی، هر توالی DNA می‌تواند به‌طور مجزا با سه توزیع

می‌باشند، به‌طور موفقیت‌آمیزی در تخمین اکسون‌ها در ژن‌ها مورد استفاده قرار می‌گیرند. با این وجود مشکل اصلی این روش‌ها این است که یک ناحیه کدکننده پروتئین ممکن است، در توالی یک ارگانیسم وجود داشته باشد؛ ولی در پایگاه داده مورد بررسی نمایش داده نشود. برای رفع این مشکل الگوریتم‌های مبتنی بر فیلتر که بر اساس تحلیل‌های طیفی [12]-[13] استوار هستند، در سال‌های اخیر به ابزارهای مفیدی در ژن‌یابی تبدیل شده است. الگوریتم‌های مختلفی برای تعیین نواحی ژنی بر مبنای خاصیت تناوب-۳ در مراجع پیشنهاد شده است [14]. از تبدیل فوریه برای این منظور استفاده کرده است. این روش با انتخاب یک پنجره با طول ثابت و لغزاندن آن بر روی توالی عددی DNA و سپس اعمال تبدیل گسسته فوریه و محاسبه انرژی طیف حاصل‌شده، نواحی ژنی را تعیین می‌کند [15]. از فیلتر معکوس میان‌گذر تیز<sup>۱</sup> با فرکانس مرکزی  $2\pi/3$  به‌منظور حذف نواحی غیرپروتئینی استفاده کرده است [16]. از یک الگوریتم جدید بر مبنای تبدیل فوریه و با به‌کارگیری پنجره بارتلت<sup>۲</sup> برای تعیین نواحی ژنی استفاده کرده است [17]. از الگوریتم‌های زمانی برای تعیین نواحی کدکننده پروتئینی استفاده کرده است [18]. از ترکیب الگوریتم گورتزل<sup>۳</sup> و تبدیل گسسته موجک<sup>۴</sup> به‌منظور تعیین نواحی ژنی و حذف نویز پس‌زمینه استفاده کرده است. در تمامی الگوریتم‌های مطرح‌شده از پنجره با طول ثابت ۳۵۱ استفاده شده است. با استفاده از پنجره با طول ثابت نواحی پروتئینی با ابعاد کوچک قابل شناسایی نمی‌باشند و از این‌رو استفاده از الگوریتم‌های مطرح با محدودیت مواجه می‌شوند.

در این مقاله روشی بر مبنای پنجره با طول متغیر جهت تعیین نواحی کدکننده پروتئین در توالی DNA ارائه شده است. استفاده از این الگوریتم منجر به افزایش دقت تخمین نواحی کدکننده پروتئین و حذف نوفه پس‌زمینه به مقدار قابل توجهی می‌شود. پیکربندی مقاله به‌صورت زیر است؛ در بخش ۲، تبدیل داده‌های نویسه‌ای DNA به سیگنال‌های دیجیتال بررسی می‌شود. در بخش ۳ تبدیل فوریه کوتاه‌مدت معرفی و مشکل آن مطرح می‌شود. نمودار الگوریتم پیشنهادی در بخش ۴ مطرح و مراحل مختلف آن توضیح داده می‌شود. ایده به‌کارگیری از فیلترهای چندطبقه به‌منظور استخراج مؤلفه تناوب-۳ در بخش ۵ توضیح داده

<sup>1</sup> Anti-Notch Filter

<sup>2</sup> Bartlett Window

<sup>3</sup> Geortzel Algorithm

<sup>4</sup> Discrete Wavelet Transform

غیر وابسته  $X_n$ ،  $Y_n$  و  $Z_n$  توصیف شود که هر یک از این سه توزیع عبارتند از [22]:

$$\begin{cases} x_i = (a_i + g_i) - (c_i + t_i) \\ y_i = (a_i + c_i) - (g_i + t_i) \\ z_i = (a_i + t_i) - (g_i + c_i) \end{cases} \quad x_i, y_i, z_i \in [-1, 1], \quad i = 1, 2, 3. \quad (1)$$

که در آن  $a_i$ ،  $c_i$ ،  $g_i$  و  $t_i$  به ترتیب، توزیع های تجمعی نوکلئوتیدهای  $A$ ،  $C$ ،  $G$  و  $T$  است. در روابط یادشده مقادیر اولیه  $a_i$ ،  $c_i$ ،  $g_i$  و  $t_i$  صفر فرض شده است  $(a_0=c_0=t_0=g_0=0)$  [22]. هر یک از این سه توزیع دارای ویژگی های بیولوژیکی خاصی هستند.  $X_n$  نماینده حضور نوکلئوتیدهای  $A$  یا  $G$  که نشان دهنده تفاوت بین بازهای پیریمیدین / پورین ( $R/Y$ ) در طول توالی DNA است، می باشد. به طور مشابه،  $Y_n$  نشان دهنده توزیع بازهای کتو / آمینو ( $M/K$ ) در طول رشته DNA و  $Z_n$  نشان دهنده توزیع پیوندهای هیدروژنی ضعیف / قوی است [23]. برای زیرتوالی تشکیل شده از باز نخست تا باز  $n$ ام از توالی DNA، وقتی بازهای پورین ( $A$  یا  $G$ ) بیش تر از بازهای پیریمیدین ( $T$  یا  $C$ ) باشد،  $X_n > 0$  و در غیر این صورت  $X_n < 0$  است. وقتی مقدار بازهای پورین و پیریمیدین مساوی باشد،  $X_n = 0$  است. به طور مشابه، وقتی بازهای آمینو ( $A$  یا  $G$ ) بیشتر از بازهای کتو ( $T$  یا  $C$ ) باشد،  $Y_n > 0$  و در غیر این صورت  $Y_n < 0$  است و اگر این بازها برابر باشند،  $Y_n = 0$  است و در نهایت وقتی بازهای پیوند هیدروژنی قوی ( $A$  یا  $T$ ) بیشتر از بازهای پیوند هیدروژنی ضعیف ( $G$  یا  $C$ ) باشد،  $Z_n > 0$  و در غیر این صورت  $Z_n < 0$  است و در صورتی که مقدار این بازها باهم برابر باشد،  $Z_n = 0$  است. بنابراین منحنی  $Z$  حاوی همه اطلاعاتی است که توالی DNA متناظر با خود دارد. تحلیل توالی DNA می تواند با استفاده از مطالعه منحنی  $Z$  متناظر انجام شود. برخی از ویژگی های منحنی  $Z$  عبارتند از [22]:

- برای یک توالی DNA، یک منحنی  $Z$  منحصر به فرد متناظر با آن وجود دارد و بالعکس.
- مختصات انتهایی یک منحنی  $Z$  برای یک توالی DNA فقط به موقعیت نوکلئوتیدها در آن بستگی دارد و مستقل از ترکیب توالی بازهای DNA است.
- اگر احتمال وقوع فرکانس های نوکلئوتیدها برای چهار باز در رشته DNA یکسان باشد، منحنی  $Z$  متناظر با آن تشکیل یک حلقه می دهد که از مبدأ شروع و به مبدأ ختم می شود.

- طول منحنی  $Z$  با تعداد بازها در توالی DNA متناسب و مستقل از ترتیب توالی یا ترکیب بازها در توالی است.

- یک ویژگی منحصر به فرد منحنی  $Z$  این است که در آن سه توزیع  $X_n$ ،  $Y_n$  و  $Z_n$  مستقل از یکدیگرند که این امر پیچیدگی محاسباتی مربوط به طیف DNA را کاهش می دهد.

- میانگین  $X_n$  و  $Y_n$  یا  $X_n$  و  $(-Y_n)$  به ترتیب اختلاف  $AT$  یا اختلاف  $GC$  نامیده می شود که میزان افزایش  $A$  در مقابل  $T$  و یا  $G$  در مقابل  $C$  در زیرتوالی تشکیل شده از باز نخست تا باز  $n$ ام توالی مورد مطالعه را نمایش می دهد.

- وقتی یک باز به توالی DNA مورد مطالعه اضافه شود، مختصات  $\alpha$ ،  $\gamma$  و  $z$  از منحنی  $z$  متناظر باید به اندازه  $+1$  یا  $-1$  افزایش یابد؛ صرف نظر از این که باز اضافه شده  $A$ ،  $C$ ،  $G$  و یا  $T$  باشد.

شکل ۵-الف، منحنی  $Z$  را برای این توالی  $F56F11.4$  نشان می دهد. همان طور که مشاهده می شود، با استفاده از منحنی  $Z$  می توان برخی از ویژگی های عمومی و محلی مربوط به موقعیت نوکلئوتیدها در یک ژنوم را به طوری که برای مشاهده کننده به صورت شهودی قابل درک باشد، نمایش داد؛ در حالی که نمایش چنین توالی هایی با طول زیاد به صورت نویسه ای و همچنین استخراج ویژگی از آنها بسیار دشوار می باشد. در شکل ۵-ب، منحنی یک بعدی حاصل از افکشن<sup>۱</sup> منحنی  $Z$  به ترتیب روی محورهای  $\alpha$ ،  $\gamma$  و  $z$  نشان داده شده است.

شکل ۶) منحنی اختلاف  $AT$  و  $GC$  را نشان می دهد. این منحنی ها به ترتیب میزان افزایش نوکلئوتید  $A$  نسبت به  $T$  و  $G$  نسبت به  $C$  را در توالی DNA نشان می دهد و برای باز نخست تا باز  $n$ ام یک توالی به ترتیب به صورت  $(x_n + y_n)/2$  و  $(x_n - y_n)/2$  تعریف می شوند. در یک تقسیم بندی کلی منحنی های اختلاف  $AT$  و  $GC$  به سه طبقه تقسیم می شوند؛ در طبقه نوع ۱ (نوع -  $AT$ )،  $A_n \approx T_n$  و  $G_n \neq C_n$ . در طبقه نوع ۲ (نوع -  $GC$ )،  $A_n \neq T_n$  و  $G_n \approx C_n$  و در طبقه نوع ۳ (نوع ترکیبی)،  $A_n \neq T_n$  و  $G_n \neq C_n$  است. همان طور که مشاهده می شود، توالی  $F56F11.4$  در طبقه نوع ۲ قرار داشته که در آن،  $A_n \neq T_n$  و  $G_n \approx C_n$  (خط آبی) و به صورت خطی (خط قرمز) کاهش می یابد.

<sup>1</sup> Projection



متمرکز بوده و سپس تحلیل‌های طیفی از جمله تبدیل فوریه کوتاه‌مدت<sup>۱</sup> به‌منظور یافتن موقعیت مؤلفه‌های متناوب بر روی توالی‌های عددی DNA اعمال می‌شود. تبدیل STFT، روشی است که در آن یک سیگنال غیرایستا<sup>۲</sup> به تعدادی از بلوک‌های کوچک‌تر تجزیه شده و سپس برای هر بلوک، تبدیل فوریه اعمال می‌شود. با استفاده از این روش، تحلیل سیگنال در حوزه زمان-فرکانس به عرض پنجره بستگی خواهد داشت؛ یعنی پنجره‌های با طول زیاد منجر به ایجاد وضوح پایین در حوزه زمان خواهد شد و بالعکس، پنجره‌های با طول کم، منجر به ایجاد وضوح بالا در حوزه زمان می‌شود. این پدیده، مهم‌ترین محدودیت تبدیل STFT به‌شمار آمده و تحلیل سیگنال‌ها را در حوزه زمان با مشکل مواجه می‌سازد. این تبدیل یک تابع چندبعدی بوده، به‌طوری‌که همواره یک سیگنال یک‌بعدی را به صفحه دوبعدی (زمان-فرکانس) نگاشت می‌دهد [24].

#### ۴- استفاده از تبدیل S در تعیین نواحی کدکننده پروتئین و الگوریتم پیشنهادی

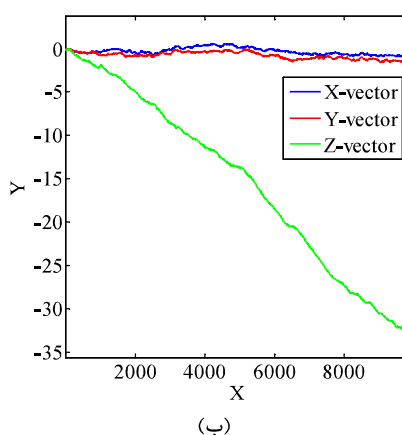
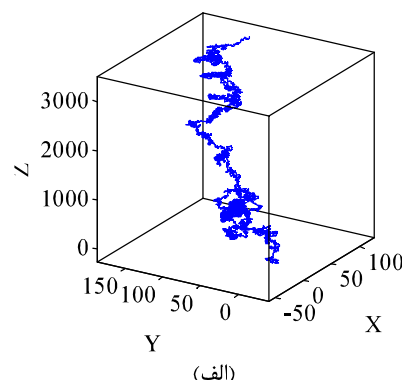
تبدیل STFT از پنجره‌هایی با طول ثابت برای محاسبات سیگنال ورودی یک‌بعدی استفاده می‌کند. در بیش‌تر مقالات این طول برابر با ۳۵۱ است. انتخاب این طول نمی‌تواند نواحی پروتئینی با ابعاد کوچک را به‌خوبی شناسایی کند. به‌منظور رفع این مشکل و بهبود عملکرد تخمین از تبدیل S استفاده می‌کنیم. این تبدیل حالت خاصی از تبدیل STFT است که در آن از یک پنجره گوسی باطول قابل تنظیم استفاده شده است. شکل (۷)، نمودار جعبه‌ای الگوریتم پیشنهادشده را نشان می‌دهد.

طبق تعریف، تبدیل فوریه کوتاه‌مدت برای یک سیگنال غیرایستان  $h(t)$  به‌صورت زیر تعریف می‌شود [25]:

$$STFT(\tau, f) = \int_{-\infty}^{\infty} h(t)g(\tau - t)e^{-j2\pi ft} dt \quad (2)$$

که در آن  $\tau$ ، موقعیت طیف زمانی،  $f$  فرکانس تبدیل فوریه و  $g(t)$  تابع پنجره گوسی<sup>۳</sup> است.

تابع گوسی قابل تنظیم به‌صورت زیر توصیف می‌شود [25]، [26]:

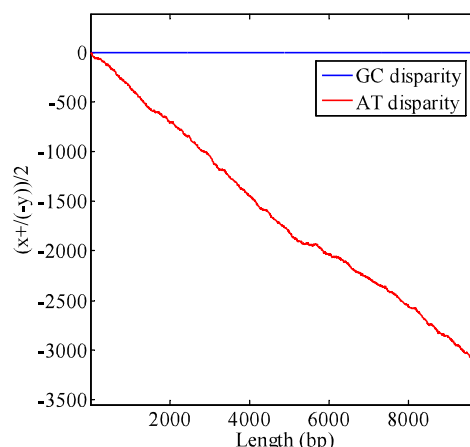


(شکل-۵): الف. منحنی سه بعدی Z برای توالی F56F11.4. ب.

منحنی یک بعدی حاصل از افکانش منحنی Z به ترتیب روی

محورهای x، y و z

(Figure-5): (a). 3-D Z-Curve for F56F11.4 sequence, (b). 1-D projected curve on x, y, and z axis, respectively.



(شکل-۶): منحنی اختلاف AT و GC در توالی F56F11.4

(Figure-6): AT and GC disparities curve for F56F11.4 sequence

#### ۳- استفاده از تبدیل فوریه کوتاه‌مدت در تعیین نواحی کدکننده پروتئین و محدودیت آن

بسیاری از روش‌های پردازش سیگنال‌های ژنومیک بر مبنای خاصیت تناوب-۳ بازهای موجود در رشته‌های DNA

<sup>1</sup> Short Time Fourier Transform

<sup>2</sup> Non-Stationary

<sup>3</sup> Gaussian Window

#### ۴-۱- خاصیت خطیگی S و تبدیل و تأثیر نوفه در آن

تبدیل S، یک عملگر خطی روی سری‌های زمانی به-شمار می‌رود. به‌طور کلی داده‌ها را می‌توان به‌صورت مجموع دو مؤلفه سیگنال و نوفه مدل کرد [28]:

$$\text{نویز} + t \text{ سیگنال} = t \quad (5)$$

با اعمال تبدیل S به رابطه (۵) داریم:

$$\{ \text{نویز} \} + S \{ \text{سیگنال} \} = S \{ \text{داده} \} \quad (6)$$

از این رو مؤلفه‌های نوفه‌ای به‌دلیل خاصیت خطیگی این تبدیل به‌طور هم‌زمان از سیگنال اصلی تفکیک شده و سیگنال و نوفه با یکدیگر ترکیب نمی‌شوند. این امر باعث افزایش دقت و بهبود عملکرد تخمین نواحی کدکننده پروتئین نسبت به روش‌های پیشین خواهد شد.

به‌طور کلی مزیت استفاده از تبدیل S بجای STFT در این است که در تبدیل S، عرض پنجره  $\sigma$  تابعی از  $f$  است. درحالی‌که در تبدیل STFT، این عرض مقدار ثابتی است. این امر باعث افزایش وضوح فرکانسی در تبدیل S می‌شود.

#### ۵- استخراج مؤلفه تناوب-۳ با به‌کارگیری فیلتر چندطبقه

به‌منظور استخراج مؤلفه تناوب-۳ از فیلتر میان‌گذر باند-محدود استفاده می‌کنیم. فیلتر میان‌گذر استفاده‌شده در این مقاله از نوع فیلترهای چندطبقه است. فیلتر چندطبقه [17] روش مؤثری در طراحی فیلتر میان‌گذر است که با افزایش تعداد ضرب‌کننده‌ها منجر به افزایش تضعیف باند توقف و بهبود عمل فیلتر می‌شود. ایده اصلی این فیلتر در شکل (۸) نشان داده شده که در آن  $H_1(z)$ ، یک فیلتر پایین‌گذر مستطیلی و  $H_1(z^3)$ ، نمونه تغییر نرخ یافته آن است. شکل (۹)، پاسخ ضربه، پاسخ دامنه و نمودار صفر-قطب فیلترهای  $H_1(z)$  و  $H_1(z^3)$  را نشان می‌دهد.  $H_1(z^3)$  باعث ایجاد دو باند عبور نامطلوب در  $\omega = 0$  و  $\omega = \frac{4\pi}{3}$  می‌شود. به‌منظور حذف این دو باند،  $H_1(z^3)$  را با فیلتر بالاگذر  $H_2(z)$  سری می‌کنیم. شکل (۱۰) پاسخ دامنه فیلتر  $H_2(z)$  را نشان می‌دهد. فیلتر  $H(z)$ ، یک فیلتر میانگذر با فرکانس مرکزی  $\frac{2\pi}{3}$  بوده که برای جداسازی مؤلفه تناوب-۳ در توالی DNA استفاده می‌شود. پاسخ دامنه این فیلتر در شکل (۱۱) نشان داده شده است.

$$g(t) = \frac{|f|}{\sigma\sqrt{2\pi}} e^{-\frac{t^2 f^2}{2}} \quad (3)$$

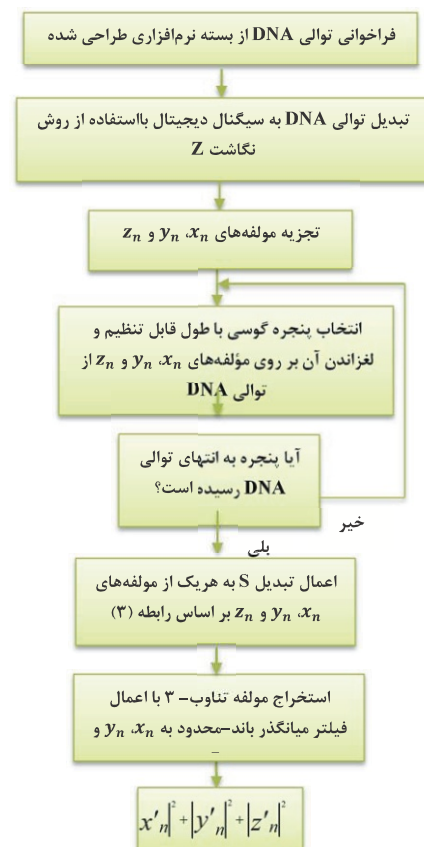
که در آن  $\sigma$  عرض پنجره است. دو ویژگی مهم پنجره گوسی استفاده شده عبارت است از [25]:

- این پنجره در هر دو حوزه زمان و فرکانس دارای تقارن است؛ به عبارت دیگر، تبدیل فوریه آن نیز یک گوسی خواهد بود.
- این پنجره فقط دارای لوب اصلی<sup>۱</sup> بوده و لوب‌های جانبی در آن وجود ندارد.

تبدیل S با جایگزینی تابع پنجره  $g(t)$  با تابع گوسی قابل تنظیم تعریف‌شده در رابطه (۳) به‌صورت زیر به‌دست می‌آید:

$$S(\tau, f) = STFT(\tau, f) = \int_{-\infty}^{\infty} h(t) \frac{f}{\sigma\sqrt{2\pi}} e^{-\frac{(\tau-t)^2 f^2}{2}} e^{-j2\pi ft} dt \quad (4)$$

مطابق رابطه (۴)، تازمانی که عرض و ارتفاع پنجره با فرکانس کنترل شود، وضوح فرکانسی در کمترین فرکانس‌ها و متناسط با آن وضوح زمانی در بیش‌ترین فرکانس‌ها به‌ترتیب با پهن‌تر و باریک‌تر شدن پنجره گوسی به‌دست می‌آید [24]، [27].



(شکل-۷): نمودار جعبه‌ای الگوریتم پیشنهادی  
(Figure-7): Block diagram of the proposed algorithm

## ۶- معیارهای ارزیابی سطح نوکلئوتیدی

به منظور مقایسه دقیق روش‌های گوناگون در تخمین نواحی کدکننده پروتئین، ارزیابی در سطح نوکلئوتیدی انجام گرفته است. در تعیین نواحی ژنی با استفاده از تکنیک‌های پردازش سیگنال، برخی از پارامترها با تغییر سطح آستانه در طیف خروجی تعریف می‌شوند. در این بخش با استفاده از شکل (۱۲)، به معرفی این پارامترها که امکان مقایسه را به دست می‌دهند، می‌پردازیم. در این شکل،  $TP$  تعداد نوکلئوتیدهایی است که به درستی به عنوان اکسون تخمین زده شده‌اند و  $TN$  تعداد نوکلئوتیدهایی است که به درستی به عنوان انترون تخمین زده شده‌اند. به طور مشابه، تعداد نوکلئوتیدهایی که به عنوان اکسون پیشگویی شده‌اند را با  $FP$  و تعداد نوکلئوتیدهایی که به عنوان انترون پیشگویی شده‌اند را با  $FN$  نشان می‌دهیم. با تعریف چهار کمیت فوق، پارامترهای حساسیت<sup>۱</sup> ( $Sn$ ) و ویژگی<sup>۲</sup> ( $Sp$ ) به صورت زیر تعریف می‌شوند [29]:

$$Sn = \frac{TP}{TP + FN} \quad (Y)$$

$$Sp = \frac{TP}{TP + FP}$$

در رابطه (Y)،  $Sn$  به صورت نسبت نوکلئوتیدهایی اکسون که به درستی به عنوان اکسون شناسایی شده‌اند و  $Sp$  به صورت نسبت نوکلئوتیدهایی اکسون تخمین زده شده که واقعاً در نواحی ژنی موجودند، تعریف می‌شود.

این پارامترها به تنهایی برای ارزیابی مناسب نیستند؛ زیرا در حساسیت بالا مقدار ویژگی پایین بوده و برعکس. بنابراین پارامتر دیگری به نام همبستگی تقریبی<sup>۳</sup> مطرح کرده که با ترکیب  $Sn$  و  $Sp$  به صورت زیر تعریف می‌شود [29]:

$$ACP = \frac{1}{4} \left( \frac{TP}{TP + FN} + \frac{TP}{TP + FP} + \frac{TN}{TN + FN} + \frac{TN}{TN + FP} \right) \quad (A)$$

$$AC = (ACP - 0.5) * 2$$

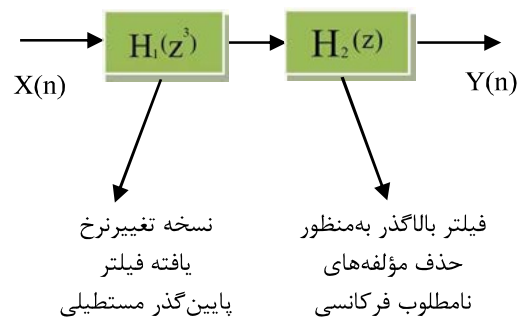
در اعمال تکنیک‌های پردازش سیگنال به منظور جستجوی نواحی ژنی، پارامترهای دیگری نیز تعریف می‌شوند. یک معیار ارزیابی متداول، منحنی ( $ROC$ )<sup>۴</sup> است که با انتخاب سطوح آستانه مختلف، مقادیر  $TP$  برای یک  $FP$  داده شده در هر سطح به دست می‌آید. ناحیه زیرمنحنی  $ROC$  نیز به عنوان یک معیار ارزیابی مورد استفاده قرار می‌گیرد؛ به طوری که مقادیر بزرگ‌تر  $AUC$  دقت بالای الگوریتم را نتیجه می‌دهد [30].

<sup>1</sup> Sensitivity

<sup>2</sup> Specificity

<sup>3</sup> Approximation Correlation

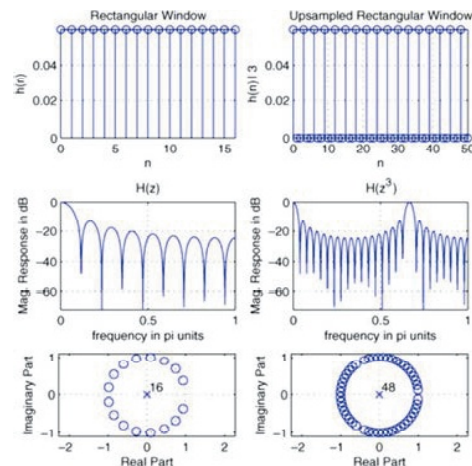
<sup>4</sup> Receiver Operating Characteristic



$$H(z) = H_1(z^3) \cdot H_2(z)$$

(شکل-۸): ایده فیلترینگ چندطبقه

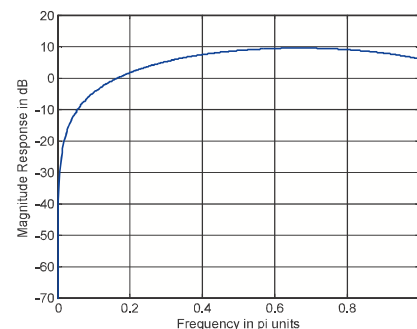
(Figure-8): Multistage filtering idea



(شکل-۹): پاسخ ضربه، پاسخ دامنه و نمودار صفر-قطب

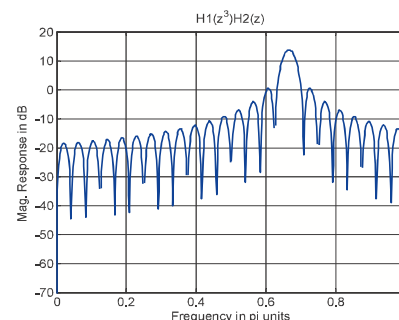
فیلترهای  $H_1(z)$  و  $H_1(z^3)$  [17]

(Figure-9): Impulse response, amplitude response, and zero-pole diagram of  $H_1(z)$  and  $H_1(z^3)$  filter [17]



(شکل-۱۰): پاسخ دامنه  $H_2(z)$

(Figure-10): Amplitude response of  $H_2(z)$



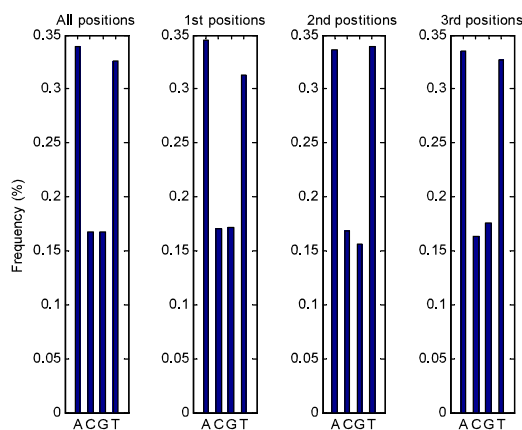
(شکل-۱۱): پاسخ دامنه فیلتر میانگذر  $H(z) = H_1(z^3) H_2(z)$

(Figure-11): Amplitude response of band-pass filter



سیگنال، امکان اعمال انواع تبدیل‌های طیفی و عددی را به منظور تحلیل توالی‌های DNA به دست می‌دهند. برخی از قابلیت‌های بسته نرم‌افزاری طراحی شده عبارتند از:

- فراخوانی انواع توالی‌های DNA در فرمت‌های مختلف و نمایش آن‌ها به صورت گرافیکی و نمادین،
- تبدیل توالی DNA به سیگنال‌های عددی با استفاده از روش‌های  $EIIP$ ، دودویی، نگاشت  $Z$ ،
- اعمال روش‌های پردازش سیگنال نظیر  $DFT$ ، تبدیل ویولت، تبدیل  $S$ ،
- مشاهده طیف رنگی انواع ژنوم‌های مختلف و ابزارهای موجود به منظور بررسی رفتار آن‌ها،
- جستجوی الگوهای خاص در توالی (نظیر کدون‌های ابتدایی و انتهایی در هر قطعه از توالی DNA)،
- پیش‌گویی مکان تقریبی نواحی ژنی با استفاده از روش‌های پردازشی و غیرپردازشی.



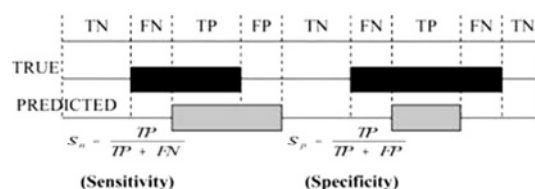
(شکل-۱۳): نمودار میله‌ای درصد حضور نوکلئوتیدها در ژن F56F11.4

(Figure-13): Bar chart of the percentage of nucleotides in F56F11.4 gene sequence

(جدول-۱): تعداد و نسبت نوکلئوتیدها در توالی ژن F56F11.4  
(Table-1): Number and proportion on nucleotides in F56F11.4 gene sequence

نوکلیوتید	تعداد نوکلئوتیدها	نسبت نوکلئوتیدها در توالی (%)
A	3332	34%
C	1647	17%
T	1647	17%
G	3206	33%

<sup>3</sup> Electron Ion Interaction Potential



(شکل-۱۲): معیارهای ارزیابی سطح نوکلئوتیدی [29]

(Figure-12): Nucleotide level evaluation criteria [29]

## ۷- نتایج پیاده‌سازی و بحث

در این بخش نتایج الگوریتم پیشنهادی و سایر روش‌های ارائه شده در مراجع با استفاده از معیارهای ارزیابی مطرح شده در بخش ۶، مقایسه شده است. این روش‌ها عبارتند از:  $AMDF$ ،  $TDP$  [17]، فیلتر معکوس میان‌گذر تیز [4]، تبدیل فوریه گسسته در زمان [14] و  $Asif$  [16].

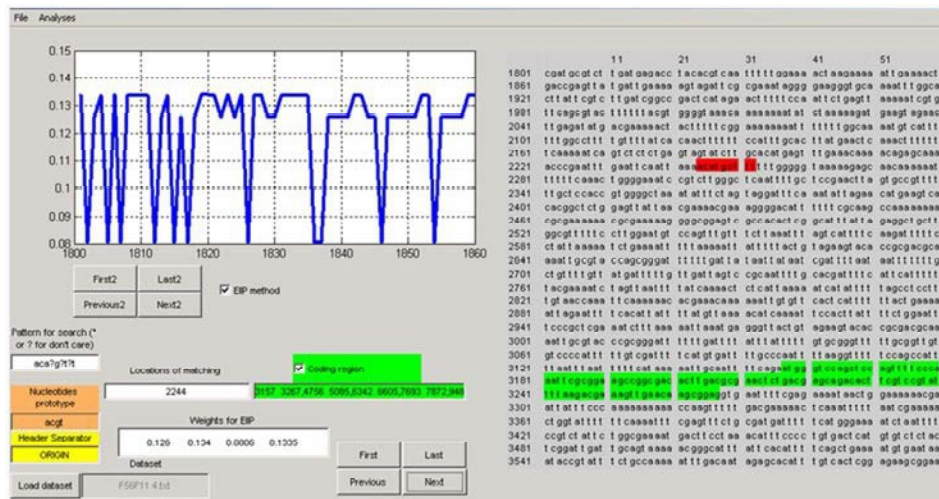
الگوریتم پیشنهادی ابتدا بر روی توالی  $F56F11.4$  در *C.elegans chromosome III* اعمال شده است. این ژن در *Caenorhabditis elegans* که یک انگل روده‌ای است که طولی در حدود ۱mm داشته و به طور طبیعی در محیط‌های معتدل خاکی زندگی می‌کند. طول بازهای آن ۱-۴۲۷۹۹ است و ۸۰۰۰ نوکلئوتید از موقعیت ۷۲۰۱ در آن وجود دارد [31].

شکل (۱۳) نمودار میله‌ای درصد حضور نوکلئوتیدها را برحسب فرکانس آن‌ها در ژن  $F56F11.4$  نشان می‌دهد. از این نمودار می‌توان به منظور محاسبه عددی مقدار محتوای  $G+C$  در نواحی پروتئینی استفاده کرد. در جدول ۱، تعداد نوکلئوتیدها و همچنین نسبت آن‌ها در توالی ژن  $F56F11.4$  نشان داده شده است. مقدار محتوای  $G+C$  در این توالی برابر ۰/۳۴ است.

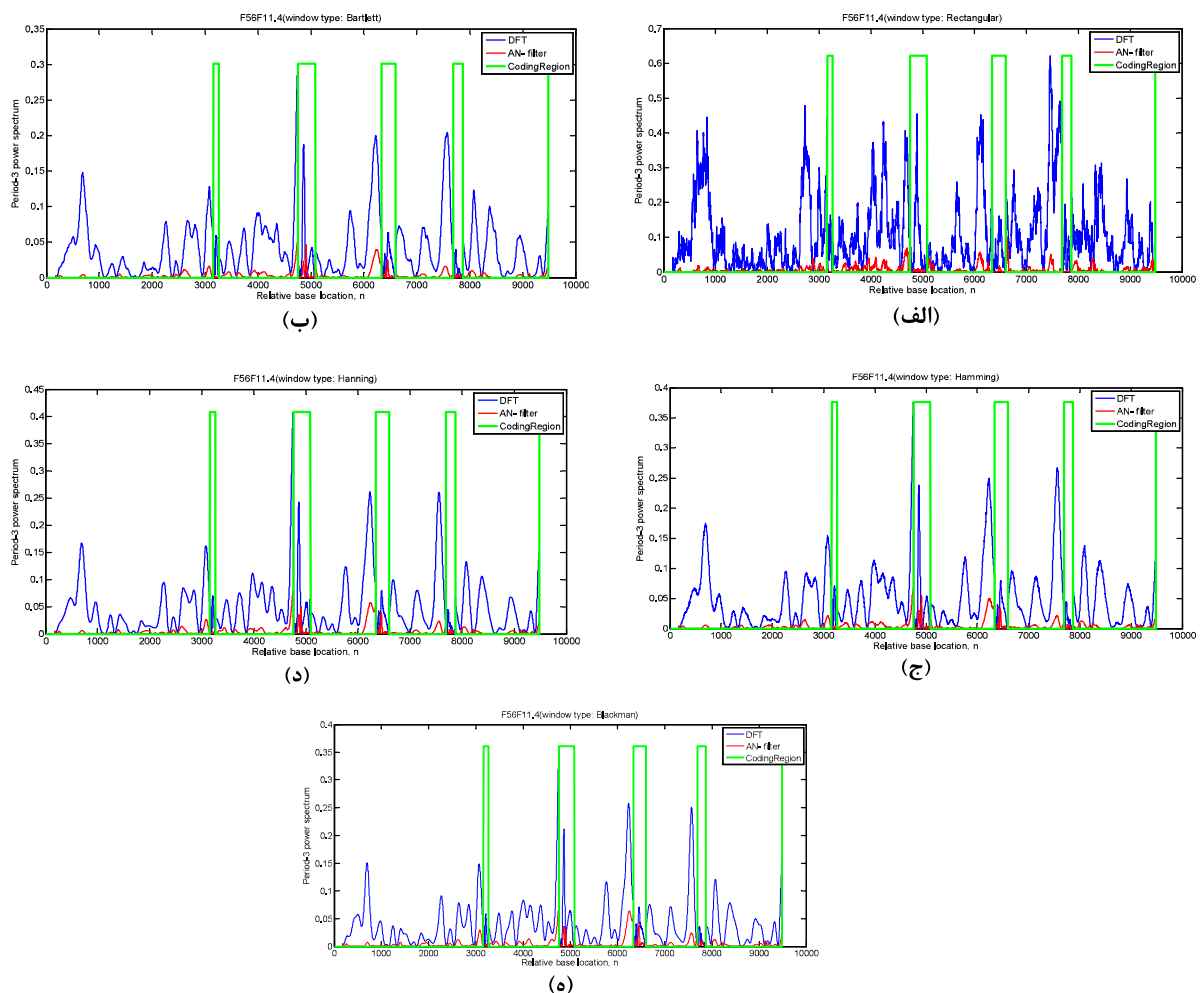
شکل (۱۴)، بسته نرم‌افزاری طراحی شده در محیط نرم‌افزار  $MATLAB$  به منظور تحلیل توالی‌های DNA را نشان می‌دهد. این بسته از دو بخش اصلی تشکیل یافته است؛ صفحه نمایش گرافیکی و ابزارهای پردازش سیگنال برای تحلیل توالی‌های DNA. صفحه نمایش گرافیکی به کاربر این اجازه را می‌دهد که ساختار توالی فراخوانده شده را هم به صورت گرافیکی و هم در قالب نمادین مشاهده کرده و الگوهای خاصی را در توالی جستجو کند. این بسته امکان فراخوانی انواع توالی‌های DNA را در فرمت‌های مختلف نظیر  $FASTA$ ،  $PHY$ ،  $TXT$  دارد. همچنین ابزارهای پردازش

<sup>1</sup> Average Magnitude Difference Function

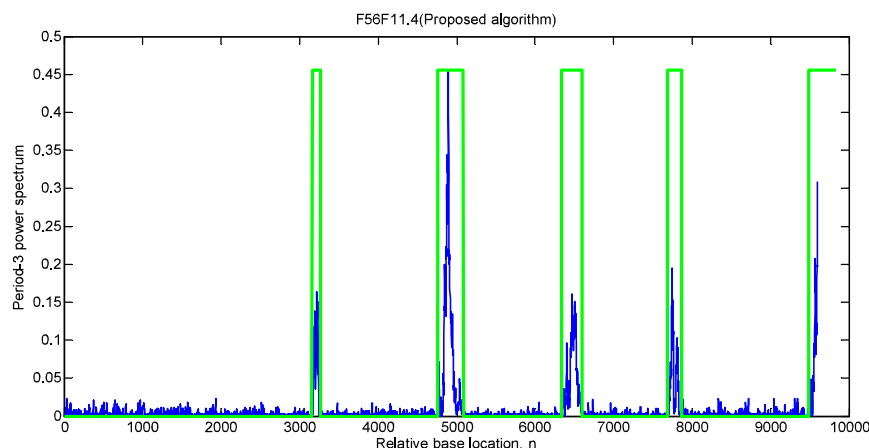
<sup>2</sup> Time Domain Periodogram



(شکل-۱۴): بسته نرم‌افزاری طراحی شده به منظور تحلیل عددی و طیفی توالی‌های DNA  
(Figure-14): Designed user friendly package in order to numerical and spectral analysis of DNA sequences



(شکل-۱۵): روش‌های DFT و فیلتر معکوس میانگذر تیز در تعیین نواحی کدکننده پروتئین با انتخاب پنجره با طول ثابت ۳۵۱ در حالات (الف): پنجره مستطیلی و (ب): پنجره بارتلت (ج): پنجره همینگ (د): پنجره هنینگ (ه): پنجره بلکمن  
(Figure-15): DFT and Anti-notch filter approaches in locating the protein coding regions by selecting the window with the length of 351 in (a). Rectangular window, (b). Bartlett window, (c). Hamming window, (d). Hanning window, and (e). Blackman window



(شکل-۱۶): اعمال الگوریتم پیشنهادی با انتخاب پنجره گوسی با طول قابل تنظیم در ژن F56F11.4  
(Figure-16): Applying the proposed algorithm by selecting the variable length window in F56F11.4 gene sequence

نشان داده شده است. برتری الگوریتم پیشنهادی نسبت به دو روش دیگر به وضوح در این جدول مشاهده می شود؛ به طوریکه در  $Sn$  برابر ۲۰ درصد، مقدار  $FP$  در الگوریتم پیشنهادی برابر ۱۲۸ جفت باز است؛ درحالی که این مقدار برای روش های  $AN$  و  $TDP$  به ترتیب برابر ۱۵۷ و ۱۹۶ جفت باز است. همچنین الگوریتم پیشنهادی مقدار  $AC$  را نسبت به روش های  $AN$  و  $TDP$  به ترتیب به میزان ۳٪ و ۶٪ بهبود می بخشد.

در نهایت الگوریتم پیشنهادی بر روی مجموعه ژن های موجود در دو پایگاه داده [32] *HMR195* و [29] *BG570* اعمال و نتایج آن با سایر روش های مطرح شده مقایسه شده است. در جدول (۵) ویژگی های این دو پایگاه داده به طور خلاصه ذکر شده است. پایگاه *HMR195* حاوی ۱۹۵ توالی ژن مربوط به انسان، موش و موش صحرایی است که در سال ۲۰۰۱ توسط *Rogic* و همکارانش تهیه شد. پایگاه *BG570* نیز حاوی ۵۷۰ توالی ژن است که توسط *Burset* و *Guigo* تهیه شده است.

(جدول-۲): مقایسه الگوریتم های مختلف در توالی F56F11.4  
(Table-2): Comparison of different algorithms in F56F11.4 gene sequence.

Method	Sn	Sp	AC
DFT	0.80	0.18	0.09
AN-filter	0.80	0.24	0.26
Asif	0.80	0.19	0.13
AMDF	0.80	0.21	0.20
TDP	0.80	0.50	0.56
Proposed	0.80	0.83	0.79

شکل های ۱۵، (الف) تا (ه) به ترتیب نتایج حاصل از روش های  $DFT$  و فیلتر معکوس میان گذر تیز را بر روی ژن *F56F11.4* با انتخاب پنجره های مختلف نشان می دهد. این پنجره ها عبارتند از: مستطیلی<sup>۱</sup>، بارتلت، همینگ<sup>۲</sup>، هنینگ<sup>۳</sup> و بلکمن<sup>۴</sup>. پنجره بلکمن به دلیل اختصاص دادن وزن زیاد به کدون های موجود در مرکز پنجره نسبت به سایر پنجره ها نتیجه بهتری به دست می دهد. در تمام حالات از پنجره به طول ۳۵۱ استفاده شده است. انتخاب این طول نمی تواند نواحی پروتئینی با ابعاد کوچک (اکسون موجود در ناحیه ۳۰۰۰ جفت باز) را به خوبی شناسایی کند. درحالی که این مشکل در الگوریتم پیشنهادی وجود ندارد؛ همانطور که در شکل (۱۶) نشان داده شده، مهم ترین ویژگی این الگوریتم، حذف نوفه پس زمینه به میزان بسیار زیاد و همچنین شناسایی نواحی کدکننده با ابعاد کوچک است.

در جدول (۲)، مقادیر  $AC$  و  $Sp$  به ازای یک مقدار ثابت  $Sn$  در الگوریتم پیشنهادی و سایر الگوریتم ها در توالی *F56F11.4* نشان داده شده است. همانطور که مشاهده می شود الگوریتم پیشنهادی بیشترین مقدار این دو پارامتر را دارا است.

الگوریتم پیشنهادی بر روی کروموزوم سوم از *Caenorhabditis* که در مجموع شامل ۱۳۷۸۳۶۸۱ نوکلئوتید با ۸۱۷۲ ناحیه کدکننده است [31]، اعمال و نتایج آن با دو روش  $AN$  و  $TDP$  مقایسه شده که در جدول (۳)

<sup>1</sup> Rectangular Window

<sup>2</sup> Hamming Window

<sup>3</sup> Hanning Window

<sup>4</sup> Blackman Window

(جدول-۳): مقایسه الگوریتم‌های مختلف در کروموزوم III از *C.elegans*  
(Table-3): Comparison of different algorithms in chromosome 3 of *C. elegans*.

Sn										
%20					%40			%60		
Methods	AUC	Fp	Sp	AC	Fp	Sp	AC	Fp	Sp	AC
AN-filter	0.6471	157	71	0.17	372	66.3	0.21	727	60	0.20
TDP	0.6115	196	70	0.15	436	65	0.18	796	59	0.19
Proposed	0.7091	128	80.5	0.24	293	76.5	0.29	615	61	0.30

(جدول-۴): خلاصه‌ای از دو پایگاه داده HMR195 و BG570  
(Table-4): Summarizing of two HMR195 and BG570 datasets

داده	ارگانیزم	# توالی ژن	جفت باز	# exon	پروتئینی	چگالی نواری
HMR195	پستانداران	195	1383720	948	14	
BG570	مهره‌داران	570	2892149	2649	15.37	

(جدول-۵): مقایسه الگوریتم پیشنهادی با سایر روش‌های موجود در دو پایگاه داده BG570 و HMR195  
(Table-5): Comparison of proposed algorithm and other different methods for two HMR195 and BG570 datasets.

Methods	BG570								HMR195					
	Sn								Sn					
	AUC	Fp	Sp	Fp	Sp	Fp	Sp	AUC	Fp	Sp	Fp	Sp	Fp	Sp
DFT	0.6540	279	45.8	767	43.3	1412	34.3	0.6782	438	51.5	1184	45	2064	41.7
AN-filter	0.6765	121	55	499	49.7	1103	36.7	0.7615	151	64.4	526	57.4	1217	51.1
Asif	0.5748	140	34.2	330	31.7	554	29.1	0.6261	214	47.1	473	44.6	787	39.9
AMDF	0.6600	340	40.8	770	39.4	1309	35.3	0.6980	410	47.9	1010	46.8	1821	43.3
TDP	0.7560	160	62	408	56	805	49.4	0.7850	262	64.8	627	60.4	1128	56
<b>Proposed</b>	<b>0.8743</b>	<b>79</b>	<b>80.5</b>	<b>214</b>	<b>72</b>	<b>534</b>	<b>64</b>	<b>0.8850</b>	<b>134</b>	<b>76</b>	<b>392</b>	<b>72</b>	<b>828</b>	<b>64</b>

(جدول-۶): مقادیر همبستگی تقریبی برای داده‌های ژنی موجود در دو پایگاه داده BG570 و HMR195  
(Table-6): Quantities of approximation correlation for gene sequences of two HMR195 and BG570 datasets

Method	BG570			HMR195		
	Sn	Sp	AC	Sn	Sp	AC
DFT	0.80	0.28	0.18	0.80	0.31	0.18
AN-filter	0.80	0.26	0.17	0.80	0.39	0.32
Asif	0.80	0.25	0.10	0.80	0.30	0.15
AMDF	0.80	0.29	0.20	0.80	0.37	0.27
TDP	0.80	0.37	0.31	0.80	0.44	0.38
<b>Proposed</b>	<b>0.80</b>	<b>0.48</b>	<b>0.45</b>	<b>0.80</b>	<b>0.51</b>	<b>0.53</b>

(جدول-۷): مقایسه کمی زمان اجرای الگوریتم پیشنهادی با سایر روش‌های موجود در پایگاه‌های داده GenBank، HMR195 و BG570  
(Table-7): Comparison of execution time of proposed algorithm and other methods for gene sequences in GenBank, HMR195 and BG570 datasets.

Gene identifier	Sequence Length (bp)	Average Computational Time (Second)		
		Proposed method	AN-filter	DFT
F56F11.4	9833	<b>695.2015</b>	714.6968	718.4017
AF009962	7422	<b>450.0254</b>	712.2368	391.0609
AJ223321.1	5321	<b>639.3497</b>	710.5081	546.2907

## ۸- نتیجه‌گیری

در این مقاله، یک الگوریتم جدید بابه‌کارگیری پنجره گوسی با طول قابل تنظیم و بر مبنای منحنی سه بعدی Z به منظور افزایش دقت در تخمین نواحی کدکننده پروتئین ارائه شد. استفاده از ایده فیلترینگ چند طبقه به منظور استخراج مولفه تناوب-۳، امکان شناسایی نواحی ژنی را آسان‌تر می‌سازد. یک مزیت مهم الگوریتم پیشنهادی حذف نوفه پس‌زمینه به میزان بسیار زیاد در آن به علت استفاده از پنجره گوسی با طول قابل تنظیم است. توانایی آشکارسازی نواحی اکسون با ابعاد کوچک نیز یکی دیگر از مزایای این الگوریتم است. با مقایسه الگوریتم پیشنهادی با سایر روش‌های موجود، مشاهده می‌شود که این الگوریتم برای داده‌های HMR195 و BG570، سطح زیر منحنی ROC را از ۴ درصد تا ۴۰ درصد بهبود می‌بخشد. روش پیشنهادی ما همچنین تعداد نوکلئوتیدهای نادرستی را که به عنوان نواحی کدکننده تخمین زده شده‌اند کاهش می‌دهد. این کاهش تعداد نوکلئوتیدهای نادرست منجر به افزایش میزان Sp خواهد شد. برای مثال، در Sn برابر ۳۰ درصد، میزان بهبود مقدار Sp در الگوریتم پیشنهادی نسبت به سایر روش‌ها از ۱۵ درصد تا ۸۵ درصد است.

نتایج حاصل شده از اعمال الگوریتم پیشنهادی و سایر روش‌ها در جدول (۵) بر روی داده HMR195 نشان داده شده است. همان‌طور که مشاهده می‌شود، الگوریتم پیشنهادی کمترین تعداد نوکلئوتیدهای نادرست شناخته شده به عنوان اکسون را دارد. در Sn برابر با ۳۰ درصد، تعداد نوکلئوتیدهای نادرست در الگوریتم پیشنهادی با ضریب ۱/۲۱ نسبت به بهترین روش مطرح شده در مراجع، یعنی روش Asif، بهبود می‌یابد. همچنین الگوریتم پیشنهادی، سطح زیر منحنی ROC را نسبت به روش‌های DFT، فیلتر AN، Asif و AMDF و TPD به ترتیب به میزان ۲۳٪، ۱۴٪، ۲۹٪، ۲۱٪ و ۱۱٪ بهبود می‌بخشد. برتری مشابهی نیز در الگوریتم پیشنهادی برای داده BG570 وجود دارد که در جدول (۵) نشان داده شده است.

جدول (۶) مقدار AC را در الگوریتم پیشنهادی و سایر روش‌های مطرح شده نشان می‌دهد. در Sn برابر با ۸۰ درصد، مقدار AC در الگوریتم پیشنهادی برابر با ۴۵ درصد در داده BG570 است؛ در حالی که مقدار آن در روش TDP برابر ۳۱ درصد است.

در جدول (۷) زمان اجرای الگوریتم و مقایسه آن با سایر روش‌ها برای توالی‌های F56F11.4، AJ223321.1 از پایگاه HMR195 و AF009962 از پایگاه BG570 نشان داده شده است. لازم به ذکر است که پیچیدگی الگوریتم پیشنهاد شده در این مقاله به دلیل استفاده از پنجره با طول متغیر، در حدود تبدیل فوریه گسسته در زمان و برابر  $O(n^2)$  است. برای این منظور در توالی ژن F56F11.4، ابتدا این توالی به پنج زیرتوالی تقسیم و بر روی هر یک الگوریتم پیشنهاد شده اعمال شده است. با انجام این کار پیچیدگی الگوریتم به  $O(n \log n)$  تنزل پیدا می‌کند.

## 9-References

## ۹-مراجع

- [1] D. P. Snustad, and M. J. Simmons, *Principles of Genetics*, John Wiley & Sons Inc, 2000.
- [2] E. R. Dougherty, et al., *Genomic signal processing and statistics*, EURASIP Book Series on Signal Processing and Communications, 2005.
- [3] D. L. Brutlag, *Understanding the human genome*, Eds. New York: Scientific American, 1994.



Applications (ISIEA), Bandung, Indonesia, pp. 354-359, September 2012.

- [16] S. Datta, A. Asif, "A Fast DFT-Based Gene Prediction Algorithm for Identification of Protein Coding Regions," Proceedings of the 30<sup>th</sup> International Conference on Acoustics, Speech, and Signal Processing 2005.
- [17] M. Akhtar, J. Epps, E. Ambikairajah, "Signal Processing in sequence Analysis: advanced in Eukaryotic gene Prediction," *IEEE journal of selected topics in signal processing*, vol. 2, pp. 310-321, 2008.
- [18] H. Saberhari, M. Shamsi, H. Heravi, and M. H. Sedaaghi, "A Fast Algorithm for Exonic Regions Prediction in DNA," *Journal of Medical Signals and Sensors*, vol. 3, no. 3, pp. 139-149, 2013.
- [19] J. M. Claverie, "Computational methods for the identification of genes in vertebrate genomic sequences," *Hum. Mol. Genet.*, vol. 6, no. 10, PP. 1735-1744, 1997.
- [20] W. F. Doolittle, "Phylogenetic classification and the universal tree," *Science*, vol. 284, no. 5423, pp.2124-2128, 1999.
- [21] F. Gao, and C. T. Zhang, "GC-Profile: a web-based tool for visualizing and analyzing the variation of GC content in genomic sequences," *Nucleic Acids Research*, vol. 34, pp. 686-691, 2006.
- [22] C. T. Zhang, and R. Zhang, "analysis of distribution of bases in the coding sequences by a diagrammatic technique," *Nucleic Acids Research*, vol. 19, pp. 6313-6317, 1991.
- [23] A. Rushdy, and J. Tuqan, "Gene Identification using the Z-Curve Representation," International Conference on Acoustics, Speech, and Signal Processing, pp. 1024-1027, 2006.
- [24] R. G. Stockwell, Why Use the S-Transform? Boulder, CO: Fields Institute Communications, 2007.
- [25] S. S. Sahu, G. Panda, "Identification of protein-coding regions in DNA sequences using a time-frequency approach," *Genomics Proteomics Bioinformatics*, vol. 9, pp. 45-55, 2011.
- [26] M. K. Hota, and V. K. Srivastava. "DSP Technique for Gene and Exon Prediction Taking EHP Indicator Sequence," In Proceedings of the Second International
- [4] P. P. Vaidyanathan, and B. J. Yoon, "The role of signal processing concepts in genomics and proteomics," *Journal of Franklin Institute, special issue on Genomics*, 2004.
- [5] E. N. Trifonov, and J. L. Sussman, "The pitch of chromatin DNA is reflected in its nucleotide sequence," *Proc. of the Nat. Acad. Sci.*, vol. 77, pp. 3816-3820, 1980.
- [6] X. F. Wan, D. Xu, A. Kleinhofs and J. Zhou, "Quantitative relationship between synonymous codon usage bias and GC composition across unicellular genomes," *BMC Evolutionary Biology*, vol. 4, no. 19, 2004.
- [7] H. E. Herzel, N. Trifonov, O. Weiss, and I. Große, "Interpreting correlations in bio-sequences," *Physica A*, vol. 249, pp. 449-459, 1998.
- [8] R. F. Voss, "Evolution of long-range fractal correlations and 1/f noise in DNA base sequences," *Phy. Re. Lett*, vol. 85, pp. 1342-1345, 1992.
- [9] C. A. Chatzidimitriou, and D. Larhammar, "Long-range correlations in DNA," *Nature*, vol. 361, pp. 212-213, 1993.
- [10] J. Henderson, J, et al., "Finding genes in DNA with a Hidden Markov Model," *J. Comput. Biol*, vol. 4, pp. 127-141, 1997.
- [11] C. H. Ding, and I. Dubchak, "Multi-class protein fold recognition using support vector machines and neural networks," *Bioinformatics*, vol. 17, pp. 349-358, 2001.
- [12] D. Anastassiou, "Genomic signal processing," *IEEE Sign. Proc. Mag*, vol. 18, pp. 8-20, 2001.
- [13] T. W. Fox, and A. Carreira, "A digital signal processing method for gene prediction with improved noise suppression," *EURASIP J. Appl. Aign. Proc*, pp. 108-114, 2004.
- [14] S. Tiwari S, S. Ramachandran, A. Bhattacharya, S. Bhattacharya, R. Ramaswamy, "Prediction of probable genes by Fourier analysis of genomic sequences," *Comput Appl Biosci*, vol. 13, pp. 263-270, 1997.
- [15] H. Saberhari, M. Shamsi, M. H. Sedaaghi, and F. Golabi, "Prediction of protein coding regions in DNA sequences using signal processing methods," 2012 IEEE Symposium on Industrial Electronics and



**موسی شمس** مدرک کارشناسی و کارشناسی ارشد خود را به ترتیب در سال های ۱۳۷۵ و ۱۳۷۹ در رشته های مهندسی برق- الکترونیک و مهندسی پزشکی از دانشگاه تبریز و دانشگاه

تهران اخذ و همچنین مدرک دکترای خود را در گرایش مهندسی پزشکی در سال ۱۳۸۷ از دانشگاه تهران دریافت کرد. زمینه های پژوهشی مورد علاقه ایشان پردازش تصویر و سیگنال حیاتی، شناسایی آماری الگو و الگوریتم های بهینه سازی بوده و در حال حاضر عضو هیأت علمی با مرتبه دانشیاری در دانشگاه صنعتی سهند است.

نشانی رایانامه ایشان عبارت است از:

shamsi@sut.ac.ir



**محمد حسین صدیقی** مدرک کارشناسی و کارشناسی ارشد خود را به ترتیب در سال های ۱۳۶۵ و ۱۳۶۶ در رشته مهندسی برق از دانشگاه صنعتی شریف اخذ و مدرک دکترای

خود را در رشته مهندسی برق در سال ۱۳۷۷ از دانشگاه لیورپول انگلستان دریافت کرد. زمینه های پژوهشی مورد علاقه ایشان پردازش تصویر و سیگنال، شناسایی آماری الگو و بیومتریک بوده و در حال حاضر عضو هیأت علمی با مرتبه استادی در دانشگاه صنعتی سهند است.

نشانی رایانامه ایشان عبارت است از:

sedaaghi@sut.ac.ir

Conference on Information Processing, 117-123. Piscataway, NJ: IEEE Press, 2008.

[27] R. G. Stockwell, L. Mansinha, and R. P. Lowe. "Localization of the Complex Spectrum: The S-Transform," *IEEE Trans. on Sig. Process.* vol. 44, no. 4, 998-1001, 1996.

[28] Yu. H. Wang, "The Tutorial: S-Transform," Accessed June 30, 2011.

[29] M. Burset, R. Guigo, "Evaluation of gene structure prediction programs," *Genomics*, pp. 353-367, 1996.

[30] M. Akhtar, E. Ambikairajah, J. Epps, "Detection of period-3 behavior in genomic sequences using singular value decomposition," *Proceedings of the International Conference on Emerging technologies*, 2005.

[31] National Center for Biotechnology Information, National Institutes of Health, National Library of Medicine, <http://www.ncbi.nlm.nih.gov/GeneBank/index.html>.

[32] S. Rogic S, A. K. Mackworth, and B. F. Ouellette, "Evaluation of gene-finding programs on mammalian sequences," *Genome*, vol. 11, pp. 817-832, 2001.

**حمیدرضا صابرکاری** تحصیلات



مقطع کارشناسی خود را در سال ۱۳۹۰ در رشته مهندسی برق- الکترونیک در دانشگاه گیلان به پایان رساند. تحصیلات کارشناسی ارشد خود را در رشته مهندسی

برق- مخابرات سیستم در دانشگاه صنعتی سهند تبریز ادامه داد و در سال ۱۳۹۲ مدرک کارشناسی ارشد خود را اخذ کرد. ایشان در حال حاضر در مقطع دکترای تخصصی رشته مهندسی برق- الکترونیک دانشگاه صنعتی سهند تبریز مشغول به تحصیل است. زمینه های پژوهشی مورد علاقه ایشان عبارتند از: پردازش سیگنال و تصویر حیاتی، بیوانفورماتیک و حسگرهای زیستی.

نشانی رایانامه ایشان عبارت است از:

h\_saberkari@sut.ac.ir