

## روش تصمیم‌گیری پس‌خور در شناسایی زنجیره‌های DNA

سارا حاجی کاظمی، حسین ضمیری و مرتضی خادمی

دانشگاه فردوسی مشهد، دانشکده مهندسی، گروه برق

E-mail: sarah.hajikazemi@gmail.com, hzamiri@ferdowsi.um.ac.ir, mortezakhademi@yahoo.com

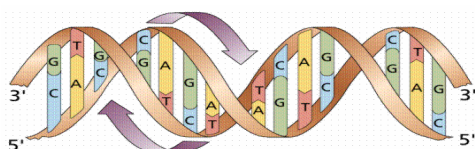
چکیده: زنجیره DNA (deoxyribonucleic acid) شامل اطلاعات ژنتیکی هر فرد می‌باشد. استخراج این اطلاعات به عنوان توالی‌یابی DNA (DNA sequencing) شناخته می‌شود. از آنجا که خطاهای کوچک در فرایند استخراج اطلاعات سبب ایجاد اشتباهات بزرگ در آزمایشات ژنتیکی می‌شود، دقت در طی فرایند توالی‌یابی برای کاهش خطاها بسیار حایز اهمیت می‌باشد برای مدلسازی این خطاها و کاهش احتمال خطا در طی این فرآیند از تئوری آشکار سازی در سیستم‌های مخابرات دیجیتال کمک گرفته می‌شود. در این مقاله روش آشکار سازی زنجیره‌های DNA با به کارگیری فیلتر منطبق تغییر پذیر با زمان و استفاده از روش تصمیم‌گیری پس‌خور ارائه شده است. شبیه سازی‌ها نشان می‌دهد که روش تصمیم‌گیری پس‌خور پیشنهادی دارای کارایی خیلی بهتری در مقایسه با روش معمولی آشکار سازی زنجیره DNA است و احتمال خطا را به طور قابل توجهی کاهش می‌دهد.

کلید واژه: DNA، توالی‌یابی DNA (DNA sequencing)، اکوالایزر تصمیم‌گیر با پس‌خور (Decision feedback equalizer)، سیستم‌های متغیر با زمان.

### ۱- مقدمه

شکل ۱ دیده می‌شود DNA به طور معمول از دو رشته (double strand) به هم پیچیده شده تشکیل شده است. هر یک از این رشته‌های تک‌تکی (single strand) از ۴ نوع باز متفاوت تشکیل شده است که بر اساس نوع اطلاعات ژنتیکی، ترتیب و تعداد آنها در رشته متفاوت است. این چهار نوع باز آدنین (A)، گوانین (G)، سیتوزین (C) و تیمین (T) می‌باشند. یک double strand از اتصال دو single strand که مکمل یکدیگرند تشکیل شده است. A و T مکمل و C و G مکمل یکدیگرند.

در بیولوژی از روش‌های مختلفی برای توالی‌یابی استفاده می‌شود که همه آنها شامل چند مرحله از جمله PCR (واکنش زنجیره‌ای پلیمرز) است.

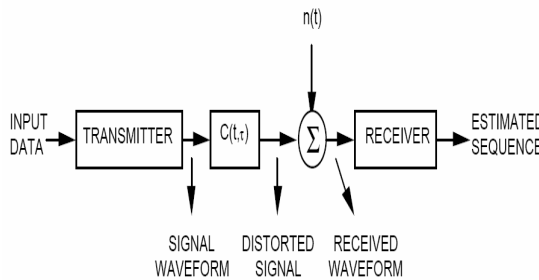


شکل ۱: ساختار مارپیچ DNA

فرآیند توالی‌یابی DNA در رأس پروژه ژنوم انسانی قرار دارد. ژنوم بیان‌کننده گروه کاملی از ساختار ژنتیکی (مجموعه‌ای از  $n$  کروموزوم) است که به وسیله یک سلول از یک ارگانیزم حمل می‌شود. دانشمندان امیدوارند که این پروژه انقلابی در علوم پزشکی و نحوه درمان بیماری‌ها ایجاد نماید. اهداف پروژه ژنوم انسانی عبارتند از تشخیص تقریبی ۲۰۰۰۰ تا ۲۵۰۰۰ ژن در ژنوم انسان، بیان توالی ۳ بیلیون جفت باز شیمیایی سازنده DNA انسان، ذخیره اطلاعات در Data base ها و بهبود روش آنالیز داده‌ها. توالی‌یابی ساختارهایی را برای هر آنچه که یک سلول انجام می‌دهد، از لحظه‌ی تولد تا مرگ شامل می‌شود. تغییرات در ساختار DNA به عنوان مثال می‌تواند شانس مبتلا شدن به بیماری را افزایش دهد و در قدرت مقابله بدن با بیماری و یا مقاومت بدن در برابر داروهای تغییراتی ایجاد نماید. با توالی‌یابی دقیق DNA، دانشمندان قادر خواهند بود عامل بیماری‌ها را شناسایی و نحوه درمان آنها را بیابند. همان‌طور که در



شکل ۲: آشکار سازی نوکلئوتیدهای فلورسنت شده با لیزر

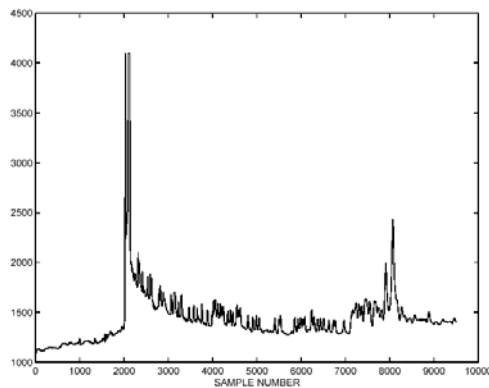


شکل ۳: شکل ۱، بلوک دیاگرام سیستم مخابرات داده

## ۲- سری زمانی DNA

در این بخش یک مدل آماری از سری زمانی DNA به دست آمده از دستگاه توالی یاب معرفی می شود. شکل پیک و پارامترها مدل شده است و نویز کل سیستم یک نویز گوسی جمع شونده فرض می شود.

به طور کلی در سرتاسر ناحیه مرکزی در شکل حاصل از دستگاه توالی یاب (شکل ۴)، یک گرایش رو به پایین در دامنه پیک ها دیده می شود. با قدرت تفکیک بیشتر برای چهار بازمتفاوت، مدل سری زمانی به صورت رابطه (۱) پیشنهاد می شود. [۱]



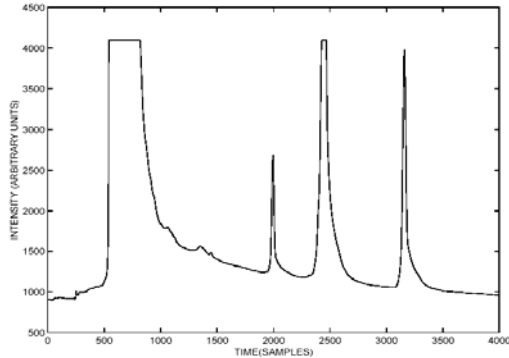
شکل ۴: شکل حاصل از دستگاه توالی یاب [۱]

در فرآیند PCR، DNA به طور نمایی زیاد می شود. مرحله بعد توالی یابی به وسیله تئوری الکتروفورسیس توضیح داده می شود که در آن چاهکهای موجود بر روی ژل به عنوان غربال هایی برای مولکولها عمل می کند به طوری که حرکت انواع مختلف اسید نوکلئیک به طول آنها بستگی دارد. به این ترتیب بر اساس تفاوت در اندازه طول، DNA ها از هم جدا می شوند. از فلورسنت رنگی برای علامتگذاری نوکلئوتیدها استفاده می کنیم که همان طور که در شکل ۲ مشاهده می شود با لیزر آشکارسازی می شوند خطاها در فرآیند توالی یابی DNA بر سه نوع می باشند.

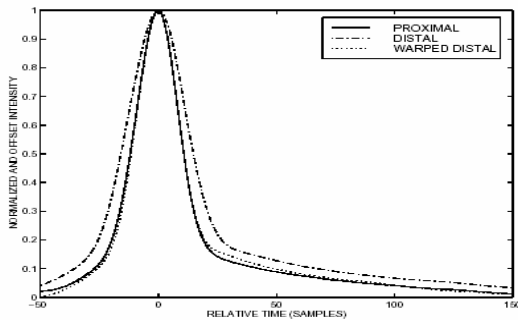
خطای جابجایی<sup>۱</sup> زمانی رخ می دهد که یک باز را به جای باز دیگر بخوانیم. خطای میان گذاری<sup>۲</sup> زمانی رخ می دهد که در اثر اشتباه یک باز را اضافه بخوانیم. خطای حذف<sup>۳</sup> زمانی رخ می دهد که در اثر اشتباه یک باز را نخوانیم. هر کدام از این خطاها می توانند بر اثر دقت پایین دستگاه توالی یاب رخ دهند. آشکارسازی توسط دستگاه توالی یابی دارای همه خطاهای فوق است. مخابرات داده بر پایه انتقال و دریافت رشته اطلاعات با احتمال کمی خطا استوار است. در شکل ۳ عناصر یک سیستم مخابرات داده نشان داده شده است. [۱] در ابتدا، داده های ورودی وارد فرستنده شده و یک سیگنال آنالوگ در ورودی کانال حاصل می شود. عبور از کانال دو اثر مهم دارد. ابتدا شکل موج توسط پاسخ ضربه کانال خراب می شود. در مرحله بعد نویز به اطلاعات اضافه می شود. گیرنده، خطاهای آشکارسازی را با تخمین زدن، حذف ISI و میانگین گیری در زمان، برای کاهش تأثیر نویز، کاهش می دهد. در این مقاله، توالی یابی DNA به صورت وابسته به نویز و رویهم افتادگی سیگنال که یک توالی اطلاعاتی را بیان می کند، مدل سازی شده است. مخابرات داده نیز برای استخراج یک توالی اطلاعاتی از حالت نویزی و روی هم افتادگی سیگنال ها به کار رفته است. به طور قطع شباهت قوی بین مخابرات داده و توالی یابی DNA وجود دارد. در بخش های بعدی ابتدا یک مدل از سری زمانی DNA و سپس سیستم مخابراتی استفاده شده ارائه می شود و نهایتاً نتایج شبیه سازی مورد بررسی قرار خواهد گرفت.

<sup>1</sup> Substitution  
<sup>2</sup> Insertion  
<sup>3</sup> Deletion

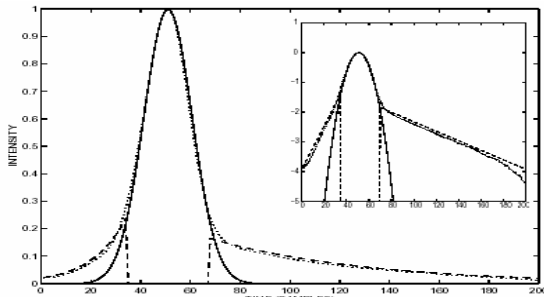
برای مدل کردن و شبیه سازی ساختار DNA از یک سیستم مخابراتی استفاده می کنیم که در بخش ۳ به توضیح آن می پردازیم.



شکل ۵: پیکهای دور از مبدا و نزدیک مبدا قبل از انطباق



شکل ۶: پیک دور از مبدا (خط پر) و نزدیک مبدا (خط نقطه)



شکل ۷: تقریب نمایی از پیک دور از مبدا

### ۳- مدل سیستم مخابراتی

همان طور که در شکل ۸ نشان داده شده است، در فرستنده به هر باز یک کد اختصاص می دهیم و بازهای DNA را با این کدها شناسایی می کنیم. A با 00، C با 01، G با 10 و T با 11 نشان داده شده است.

در محل هر باز پیک ری می دهد که تنها تفاوت پیک ها در عرض پالس و دامنه است. خطاهای ناشی از میان گذاری و حذف ناشی از نویز جمع شونده ای است که آن را سفید و گوسی مدل می کنند. همان طور که در شکل ۹ نشان داده

$$y_{n,k} = \sum_{i=1}^{N_b} a_i g_{k,t_i} \delta_{n,x_i} + n_k \quad (1)$$

که n نوع باز را نشان می دهد. اندیس k شماره نمونه و جمع روی موقعیت توالی باز i ام و  $N_b$  تعداد کل بازهایی است که توالی یابی می شوند.  $\delta_{n,x_i}$  در تشخیص این که آیا  $x_i$  یعنی باز در موقعیت i ام توالی، همان نوع کانال n است، استفاده می شود. به این ترتیب که، در صورتی که n و  $x_i$  برابر باشند، یک ودر غیر این صورت صفر است. متغیرهای  $t_i$  و  $a_i$  لرزش زمانی و نوسان دامنه را مدل می کنند و در نهایت یک نویز جمع شونده  $\{n_k\}$  داریم که نوسان سطح زمینه را معرفی می کند. الکتروفورسز یک مولکول خالص وقتی که مولکول های DNA به اندازه کافی دور از زمان بارگذاری حرکت کند، باید منجر به یک پیک گوسی شود. اگرچه در نتایج مشاهده شده از دستگاه توالی یابی، شکل پیک بسیار پیچیده تر از گوسی است. در شکل ۵ چهار نوع پیک نمونه نشان داده شده است. [۱]

در شکل ۶ خط پیک هر ناحیه حرکت کرده تا پیک دور از مبدا با پیک نزدیک مبدا تطبیق یابد.

تنها تفاوت بارز در پیک ها، بسط افتن پیک های انتهایی در زمان است. پیک ها می توانند با فرمول (۲) مدل شوند.

$$g_{k,t_i} = g_I \left( \frac{k - t_i}{P_w(t_i)} \right) \quad (2)$$

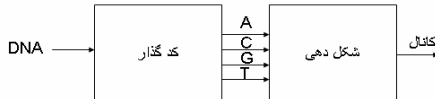
k اندیس نمونه،  $t_i$  زمان پیک،  $g_I$  شکل پالس (پیوسته) وقتی عرض پیک واحد باشد و  $P_w(t_i)$  نشان دهنده وابستگی عرض پیک به زمان پیک است که به صورت زیر مدل می شود [۱]:

$$P_w(t_i) = 15.08 + 0.0326(t_i - 1) \quad (3)$$

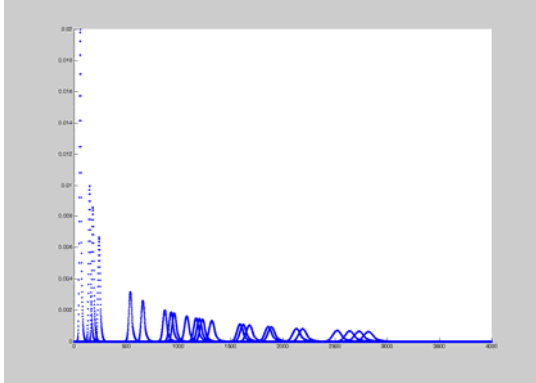
در شکل ۷، شکل پالس به سه ناحیه مجزا تقسیم می شود که لوب اصلی گوسی و دو دامنه نمایی است. [۱]

شکل پالس کلی با عرض واحد برای مجموعه داده توالی یابی، فرمول (۴) را نتیجه می دهد.

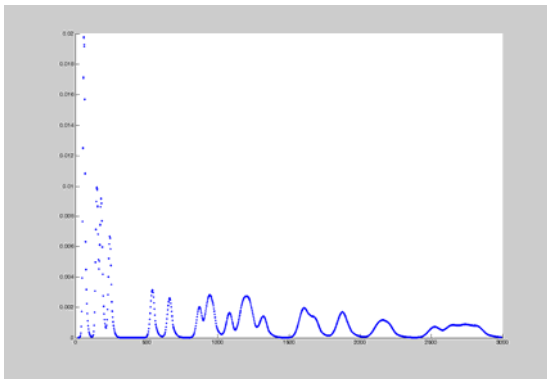
$$g_I(l) = \begin{cases} 7.89e^{4.8l}, & l < -0.86 \\ e^{-(1.67l)^2}, & -0.86 < l < 0.64 \\ 6.89e^{-1.96l}, & 0.64 < l \end{cases} \quad (4)$$



شکل ۸: کدگذاری بازها



شکل ۹: تأثیر کانال بر روی پالس های ورودی برای ۵۰ باز



شکل ۱۰: مجموع ورودی و نویز برای ۱۰۰ باز با SNR=30dB

$$P_e = \frac{P_1}{2} \operatorname{erfc} \left( \frac{\sqrt{E_k} - \lambda(k)}{\sqrt{N_0}} \right) + \frac{P_0}{2} \operatorname{erfc} \left( \frac{\lambda(k)}{\sqrt{N_0}} \right) \quad (7)$$

که در آن  $E$  تغییر پذیر با زمان بوده و آن را به صورت  $E_k$  در نظر می گیریم که از رابطه زیر حاصل می شود:

$$E_k = \int_{-\infty}^{+\infty} g_{t,k}^2 dt \quad (8)$$

چنانچه  $p_1$  احتمال وجود باز و  $p_0$  احتمال عدم وجود باز باشد برای مینیمم سازی خطا می توان نشان داد که:

$$\lambda_{opt}(k) = \frac{N_0}{2\sqrt{E_k}} \left( \ln \left( \frac{P_0}{P_1} \right) + \frac{E_k}{N_0} \right) \quad (9)$$

هر گاه  $p_0 = p_1$  باشد، داریم:

$$\lambda_{opt}(k) = \frac{\sqrt{E_k}}{2} \quad (10)$$

شده است، مطابق رابطه (۳) با افزایش تعداد نمونه ها، پالسها در حوزه زمان گسترده شده و دامنه آنها نیز کاهش می یابد. در این شکل تداخل سیگنال ها و تأثیر هر سیگنال بر روی نمونه بعدی به خوبی مشهود است.

بدیهی است که پس از جمع پالس ها در هر زمان، در هر نقطه سطح انرژی افزایش می یابد. در شکل ۱۰ دیده می شود که در نمونه های انتهایی رشته ISI ایجاد شده به همراه نویزی که در کانال به پالس ها افزوده می شود، آشکار سازی را بسیار مشکل می سازد. با توجه به دوره هر پالس در این جا می توان تخمین زد که هر پالس با دو نمونه پیش از خود تداخل می کند. برای رفع این مشکل و کاهش احتمال خطا از اکوالایزر تصمیم گیرنده با پسخور استفاده می کنیم.

شکل ۱۱ بلوک دیاگرام یک سیستم آشکارساز به روش پسخور را نشان می دهد. سیگنال  $r(t)$  سیگنال دریافتی است که جهت آشکارسازی زنجیره DNA به کار می رود. شکل ۱۰ یک نمونه نوعی از سیگنال ورودی  $r(t)$  را نشان می دهد که ناشی از وجود بازهای نشان داده شده در شکل ۹ می باشد.

همان طور که شکل ۱۱ نشان می دهد پهنای هر پالس که ناشی از وجود یک باز است، می تواند در بازه مجاور گسترش یابد. بنابراین ISI ناشی از دو باز در دو طرف باز باعث خطا در آشکارسازی می شود. نظر به این که تصمیم گیر با پسخور فقط توانایی حذف ISI ناشی از گذشته را دارد بنابراین بصورت زیر مدل سازی می گردد:

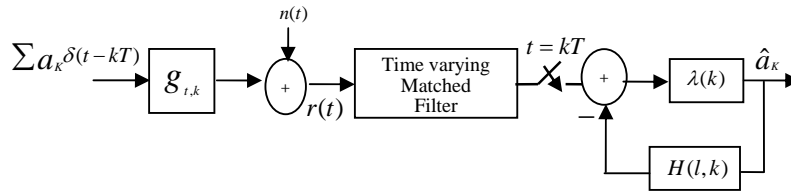
$$H(l,k) = h_{1,k} q^{-l} + h_{2,k} q^{-2} \quad (5)$$

همان طور که در رابطه (۵) دیده می شود ضرایب فیلتر پسخور وابسته به آشکارسازی باز  $k$  ام است که این ضرایب از رابطه زیر به دست می آیند.

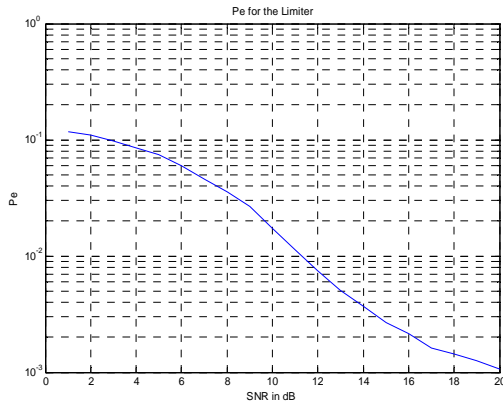
$$h_{1,k} = \int_{-\infty}^{+\infty} g(t,k) g(t,k-1) dt \quad (6)$$

$$h_{2,k} = \int_{-\infty}^{+\infty} g(t,k) g(t,k-2) dt$$

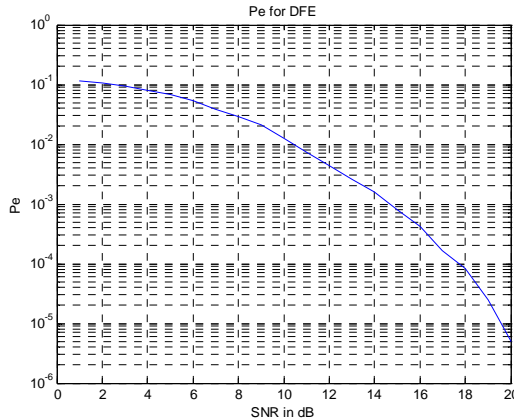
از طرفی نظر به این که انرژی  $g_{t,k}$  برای  $k$  های گوناگون متفاوت است،  $\lambda(k)$ ، سطح آستانه تصمیم گیری نیز متغیر با زمان است. هر گاه  $\lambda(k)$  سطح آستانه تصمیم گیری باشد احتمال خطا برابر است با:



شکل ۱۱: بلوک دیاگرام سیستم آشکارساز با تصمیم گیری پسخور



شکل ۱۲: تابع احتمال خطا (probability of error, Pe) برای ۱۰۰ باز بدون استفاده از روش آشکارسازی پسخور



شکل ۱۳: تابع احتمال خطا (probability of error, Pe) برای ۱۰۰ باز با استفاده از روش پسخور پیشنهادی

## ۶- مراجع

- [1] Simon Haykin, Communication Systems, 1978.
- [2] Stephen William Davis, Application of Communication Theory to Automatic DNA Sequencing, PhD thesis, 1999.
- [3] Rodney Andrew Kennedy, Operational Aspects of Decision feedback equalizers, B.E. (Hons.) UNSW, 1988.

## ۴- نتایج شبیه سازی

در شبیه سازی ها ابتدا رشته ورودی DNA را به صورت 0 و 1 مدل می کنیم. چون توالی خاصی از DNA مد نظر نیست برای حفظ کلیت، ورودی را یک مجموعه از 0 و 1 های تصادفی که دارای توزیع یکنواخت می باشد در نظر گرفته می گیریم.

در این جا بدون از دست دادن کلیت، شناسایی باز آدنین را در توالی مدنظر داریم. به طوری که به ازای هر ورودی 00 یک پیک در محل شماره باز آدنین خواهیم داشت.

شکل های ۱۲ و ۱۳ به ترتیب احتمال خطای آشکارسازی زنجیره DNA برای آشکارسازی معمولی (بدون استفاده از DFE) و با استفاده از روش تصمیم گیری پسخور را نشان می دهد. همان طور که ملاحظه می گردد احتمال خطای آشکارسازی پیشنهادی با پسخور به طور قابل توجهی کمتر از روش معمولی آشکارسازی زنجیره DNA می باشد.

## ۵- نتیجه گیری

در این مقاله با استفاده از تئوری آشکارسازی در سیستم های مخابرات داده جهت آشکارسازی زنجیره DNA بر پایه گیری با پسخور ارائه گشت. نظر به این که دامنه پالس های ناشی از وجود بازها در زنجیره DNA با افزایش طول DNA کاهش می یابد. آشکارسازی پیشنهادی با استفاده از یک فیلتر منطبق تغییر پذیر با زمان و نیز یک سیستم تصمیم گیر که سطح آستانه آن نیز تغییر پذیر با زمان است و برای مینیمم سازی احتمال خطا طراحی گشته است. صورت می پذیرد. فیلتر FIR تغییر پذیر با زمان که در حلقه تصمیم گیری پسخور قرار دارد، باعث کاهش ISI میگردد. شبیه سازی ها نشان داد که احتمال خطای آشکارسازی با روش پیشنهادی با تصمیم گیری پسخور به طور قابل توجهی نسبت به روش معمولی آشکارسازی کاهش می یابد.

[4]Claes Tidestav, The Multivariable Decision Feedback Equalizer Multiuser Detection and Interference Rejection, Uppsala University Signals and Systems Dept,1999.

[5][http://www.ornl.gov/sci/techresources/Human\\_Genome/home.shtml](http://www.ornl.gov/sci/techresources/Human_Genome/home.shtml)

[6][http://www.blc.arizona.edu/Molecular\\_Graphics/DNA\\_Structure/DNA\\_Tutorial.HTML#Components](http://www.blc.arizona.edu/Molecular_Graphics/DNA_Structure/DNA_Tutorial.HTML#Components)

[7]<http://seqcore.brcf.med.umich.edu/doc/educ/dnapr/pg1.html>

[8]<http://www.bergen.org/AAST/Projects/Gel/intro.htm>

[9]<http://content.febsjournal.org/cgi/content/full/269/15/3632#F>

[10]<http://arbl.cvmb.colostate.edu/hbooks/genetics/biotech/gels/principles.html>