

Bioinformatic Evaluation of a Sequence for Custom TaqMan® Gene Expression Assays

Overview

The Custom TaqMan® Gene Expression Assays (formerly named TaqMan® Assays-by-Design® Gene Expression Service) are custom assays that are designed, synthesized, formulated, and delivered as analytically quality-controlled primer and probe sets for gene expression assays based on sequence information submitted by the customer. The goal of this tutorial is to help the researcher evaluate the quality of their sequence information before submitting an order for a Custom TaqMan® Gene Expression Assay. Specific information is given on how to assess a sequence using a variety of on-line tools.

The [Custom TaqMan® Gene Expression Assays](#) provide the researcher the opportunity to design an assay that is not currently available through the [TaqMan® Gene Expression Assays](#) (formerly named TaqMan® Assays-on-Demand™ Gene Expression Assays) offerings. Studies that involve viral detection, species other than *H. sapiens*, *M. musculus*, *R. norvegicus*, *A. thaliana*, *D. melanogaster*, or *C. elegans*, or detection of specific pathogens are some examples of applications that would benefit from this custom design line of products. For human, mouse, rat, *Drosophila*, *C. elegans*, and *Arabidopsis* gene expression assays, the TaqMan® Assays should be used. If a particular gene target is currently not available then one should consider a custom design.

Note: Additional TaqMan® Assays are regularly added to the web site for ordering. Please visit the [Applied Biosystems web site](#) for regular updates.

Process Overview

Ordering Custom TaqMan® Assays involves the following procedures:

1. Selecting a target sequence
- 2. Assessing the quality of the sequence**
3. Preparing the submission file using the [File Builder software](#)
4. Formatting the sequence for submission
5. Submitting the order via the File Builder software or e-mail.

Step two, **Assessing the quality of the sequence**, will be covered in this tutorial.

Step 1, and 3-5: Selecting a target sequence, Preparing the submission file, Formatting the sequence for submission, and Submitting the order are covered in:

- [Custom TaqMan® Genomic Assays: Protocol: Submission Guidelines](#)
- [Online Ordering Procedures Using the File Builder Software: Quick Reference Card](#)
- [TaqMan® Assays-by-Design Service for Gene Expression Assays Quick Reference Card](#).

Assessing the Quality of the Sequence

Overview

The most important factor in the success of the Custom TaqMan[®] Gene Expression Assays is the quality of the sequence data that you submit for the design process. Sequence analysis gives one a tool to eliminate poor sequence quality so it does not adversely impact the assay. Following this section, a variety of on-line tools are presented to help assess your sequence. Consider the following when selecting your target sequence:

- ❖ [Biological significance](#)
- ❖ [Sequence length](#)
- ❖ [Sequence quality](#)
- ❖ [Masking sequences](#)
- ❖ [Uniqueness of sequence](#)

Biological Significance

When choosing sequences to submit, first consider the biological significance of the desired assay. The quality assurance on assays carried out during manufacture of the primer and probe can ensure only that the yield and content of the primers and probe meet specifications. Applied Biosystems is unable to guarantee the biological performance of the assays.

Examples:

- If you know that your gene of interest has more than one transcript (splice variants) make sure you are submitting a sequence that will detect one or all of the variants you wish to detect. On the contrary, if you only want to detect one out of five splice variants for a particular transcript, make sure that you have selected your coordinates (see *Note* below) appropriately, and masked any unwanted regions of that transcript to ensure that the assay you receive is specific only for your transcript of interest.
- If you are studying a gene that has regions of high homology to other members within a gene family, or to closely related genes, you will want to ensure specificity by using areas of sequence unique to the gene of interest and masking homologous regions with Ns.

Note: If you are studying the gene expression of a multi-exon gene, it is important to know the location of the exon junctions within the cDNA sequence that you submit for assay design. The ideal assay design is to have the TaqMan[®] MGB probe designed across an exon-exon boundary. The exon boundary location is used as a coordinate in the sequence submission process for a gene expression assay.

Sequence Length

To optimize your assay design, follow these guidelines:

- Submit a sequence length of approximately 600 bases.
Increasing the sequence length increases the assay design possibilities.
- Select the sequence so that the target site is toward the center of the submitted sequence.

Note: *Sequence length can range from 61 to 5000 bases. Short (fewer than 300 bases) sequences limit the potential number of assays that can be designed.*

Sequence Quality

To Assess the Quality of the Sequence:

1. Obtain confidence in the sequence accuracy. You want to have the most accurate sequence of your desired target before you submit the sequence to have an assay designed. Inaccurate sequences can lead to failed assays due to poor binding, or no binding, of primers or probes.

Note: *If you performed the sequencing yourself, it is strongly recommended that you perform multiple sequencing reactions to remove any ambiguities.*

2. Use other resources, such as public databases with curated sequences such as [RefSeq](#) (which contains mRNA sequences) or [dbSNP](#) (which contains documented SNPs) to determine the quality of your sequence.

Masking Sequences

The Custom TaqMan[®] Assays proprietary software for designing primers and probes will not design probes or primers to a region of sequence containing Ns. You can annotate your sequences with Ns to avoid specific regions of sequence in design (e.g., ambiguous sequences, repetitive sequences, or SNP sites), albeit the use of Ns may limit assay design.

To mask sequences:

1. You may substitute each ambiguous base with an N.

For example:

The **bolded** bases in this sequence are ambiguous:

ACGTGACGTGACGTGACGTGACGTGGATYGTGR**SR**STCCT

Where Y= C or T, R=A or G, and S= G or C; they would be substituted as:

ACGTGACGTGACGTGACGTGACGTGGAT**NGT**G**NNN**NTCCT.

2. Minimize the substitution of Ns in the sequence.

Because the Custom TaqMan[®] Assays proprietary software does not include Ns in the probe or primer, having a sequence with Ns greatly reduces the number of available primers and probes from which to select an optimal assay.

3. Ensure that Ns are not too close to the target site.

Important! No probes can be designed if Ns are too close to the target site. When designing gene expression assays, make sure that no Ns are within five bases of the target site.

Uniqueness of Sequence

After you have selected a sequence, check whether unique primers and probes can be generated for the cDNA sequence by verifying that the target sequence is unique within the organism you are studying.

1. Substitute Ns to mask small regions of repeats and SNPs. Run the sequence through a program such as [Repeat Masker](#) to detect common repetitive elements.

2. Perform a [BLAST](#)[®] search against public databases to detect regions within your sequence that have similarity to other published sequences. If there are large regions of similarity with other sequences in a gene family, use a different area of sequence that is unique to your gene of interest.

3. For gene expression assays, choose an exon-exon boundary that is unique for the transcript(s) of interest.

For Custom TaqMan® Gene Expression Assays, the TaqMan® MGB probe, when possible, should be designed across an exon-exon boundary in order to exclude the detection of genomic DNA. The exon boundaries are what will preferably serve as your coordinate(s) in your submission file. If you are working with a gene sequence that is in a public database, there are web resources* available to find exon information. One can search the nucleotide database using [Entrez](#) at NCBI or use [Vertebrate Genome Annotation](#) (VEGA), which is part of the [Ensembl](#) project.

TOOLS

I. Repeat Masker

While the use of Ns limits assay design (see [Masking Sequences](#)), it allows you to eliminate possible assay design in areas of similarity to other unrelated sequences or to regions of low complexity DNA. Neither repeat elements nor low complexity DNA should be used as potential PCR primer or probe sites since they could produce non-specific amplification or probe binding.

On average, close to 50% of the human genomic DNA sequence will be masked by RepeatMasker. It is a program that screens DNA sequences for interspersed repeats and low complexity DNA sequences. The output is a detailed annotation of the repeats that are present in the query sequence as well as a modified version of the query sequence in which all the annotated repeats have been masked (default: replaced by Ns). The masked sequence can be used for submission and can also be used in BLAST® searches.

*Examples of web sites that host RepeatMasker are:

<http://www.repeatmasker.org>

This website has a lot of useful information on the RepeatMasker program, including FAQs and documentation such as Interpreting Results, Sensitivity, and RepeatMasker uses. “RepeatMasker is most commonly used to avoid spurious matches in database searches. Generally this step is strongly recommended before doing BLASTN or BLASTX equivalent searches with mammalian DNA sequence.”

<http://woody.embl-heidelberg.de/repeatmask>

This site is a mirror of the University of Washington site above. The [repeatmask help](#) on this site has similar information to that of the University of Washington.

How to use RepeatMasker

A. Submitting your sequence / Starting your query

- You may enter your sequence by either copying and pasting your sequence into the box provided, or uploading it from a file.
- Sequences can be submitted one at a time or in batch form.
- Sequence submissions must be in [FASTA format](#) (see input format).
- When selecting ‘return format’ and ‘return method’, if you choose “html” for both, your results will be displayed in your web browser window.

- Make sure you choose the appropriate source of your DNA. The default genome library is human. Because interspersed repeats are specific to a (group of) species, it is important to select the appropriate repeat library to search.
- Click on 'Submit Sequence'.

RepeatMasker Submission

Basic Options

Large sequences will be queued, and may take a while to process.

Enter the file to process:

Or paste the sequence(s) in FASTA format:

Enter sequence(s) here

```
>BC032413.1 Homo sapiens B lymphoid tyrosine kinase, mRNA
CACCTCTGTCTGCTGCCGGCAGAAAGCCACAAGCCATGAAAAGTATTGAGATGA
GAAGAATTCATCTGGGACTGGCTTTTGGCTTTAGGATGGTGTGGAAAGTTGCTCGTT
GTCGCTAGGAGCCTGCTCCACTGTAAGGGTGTGGGATCTGAAGAGCTATGGTG
AAACACCACTGAAGCATTGCCAAGGATGGGGCTGGTAAGTAGCAAAAAGCCGGA
CAAGCAAAAAGCCATCAAAAGACAAGGACAAGCCCAATGCGAGCCCTCAAGCT
```

Select return format: html tar file links

Select return method: html email

Advanced Options

Speed/Sensitivity: rush quick default slow

<u>DNA source</u> :	Human
<u>Contamination</u> :	Rodent
	Mouse
	Rat
<u>Repeat Option</u> :	Artiodactyls and whales
	Cow
	Pig
<u>Artifact Check</u> :	Carnivore
	Cat
	Dog
<u>Alignment Op</u> :	Chicken
	Xenopus (African clawed frog)

B. Viewing your Results

- RepeatMasker returns the submitted sequence(s) with all recognized interspersed or simple repeats masked. In the masked areas, each base is replaced with an N, so that the returned sequence is the same length as the original.
- A table annotating the masked sequences as well as a table summarizing the repeat content of the query sequence will be returned to your screen. In the "html" return format all output is returned to your screen in one file.
- The masked sequence can be copied directly from the web browser.
- We strongly recommend that when any sequence is submitted for a Custom TaqMan[®] Assay, the sequence be masked for repeat elements. This will reduce the possibility of poor sequence quality impacting assays.

RepeatMasker Output

Repeat Annotations:

SW score	perc div.	perc del.	perc ins.	query sequence	position begin	position end	in query (left)	matching repeat	repeat class/family	position begin	position end	in repeat (left)	ID
216	30.8	14.3	0.0	BC032413.1	1882	1972	(279) +	MIR	SINE/MIR	84	187	(75)	1
477	0.0	0.0	0.0	BC032413.1	2199	2251	(0) +	{A}n	Simple_repeat	1	53	(0)	2

*Masked Sequence:

```
>BC032413.1 Homo sapiens B lymphoid tyrosine kinase, mRNA
CACCTCTGTCTGCTGCCGGCAGAAAAGCCACAAGCCATGAAAACCTGATTGA
GATGAGAAGAATTTCATCTGGGACTGGCTTTTGGCTTTAGGATGGTGTGGA
AGTTGCTCGTTGTCGCTAGGAGCCTGCTCCACTGTAAGGGTGTCTGGGATC
GTGCTGGCGCAGCCGGCCCGAGGAGCGGCCACCTTCGAGTTCCTGCAGT
CGGTGCTGGAGGACTTCTACACGGCCACCGAGCGGCAGTACGAGCTGCAG
CCCTAGCCGGCCGCGCCCGCCTGCGCCCGTGGCCACCTCTGCGCGGACG
ACCCCGACTTCCGTGCCATCCCAGACGGGCCGCGAAGCGGGGTGTCCGC
TGTGCCCTTTTCTCAGACCCGGAAATCCAGTGGGCAGAGGCAGCTTCGCAG
GGGGTCCCGGACGGACTCCTTCCACCGACNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
GACTGTCATCAAAGTAAGGCCCGCGTGGGCACCCCGCGTGGCCGCGC
GTCCCGCGCTCTGCGCCCTGCGTGGACCCCGCCCTGCCCGCTACAGAA
CCAGACTGGGTCCCGCGGACGCCAGCAGGGGCACCCCGAGCCTAGGCTGC
GCTCCAGCAGTGGGGCTTTTCTGCAATAAAGTCACGAGCGTTCGNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
```

Any repeat regions are automatically converted to Ns in the submitted sequence.

* Sequence was shortened for display purposes.

Summary:

```
=====
file name: RM2sequpload_2043
sequences: 1
total length: 2249 bp (2249 bp excl N-runs)
GC level: 0.00 %
bases masked: 144 bp ( 6.40 %)
```

Number & Percentage of bases masked

Repeat Elements

	number of elements*	length occupied	percentage of sequence
SINEs:			
ALUs	0	0 bp	0.00 %
<u>MIRs</u>	1	<u>91 bp</u>	4.05 %
LINEs:			
LINE1	0	0 bp	0.00 %
LINE2	0	0 bp	0.00 %
LTR elements:			
MaLRs	0	0 bp	0.00 %
ERV1	0	0 bp	0.00 %
ERV_classI	0	0 bp	0.00 %
ERV_classII	0	0 bp	0.00 %
DNA elements:			
MER1_type	0	0 bp	0.00 %
MER2_type	0	0 bp	0.00 %
Unclassified:	0	0 bp	0.00 %
Total interspersed repeats:		91 bp	4.05 %
Small RNA:	0	0 bp	0.00 %
Satellites:	0	0 bp	0.00 %
Simple repeats:	1	53 bp	2.36 %
Low complexity:	0	0 bp	0.00 %

In this example there is a stretch of sequence that is comprised of 91 bases of MIR sequence, a common repeat element. If a TaqMan® primer or probe were designed across this MIR sequence (because it was not masked before submission) the oligo could bind to any MIR sequence in the genome. This assay would not be very discriminating or specific because of the number of sequences to which it could potentially bind.

* most repeats fragmented by insertions or deletions have been counted as one element

II. **BLAST**[®] (Basic Local Alignment Search Tool)

After you have selected a target, there are several things that must be considered before submitting a sequence for a Custom TaqMan[®] Assay. Whether you have sequenced your target or taken the sequence from a sequence database, it is important to determine whether unique primers and probes can be generated for the sequence. Homologs in gene families can present a problem, as can orthologous sequences when working in a transgenic system. It is also important to identify any polymorphisms in your sequence of interest. All of these possibilities should be considered before submitting a sequence for a Custom TaqMan[®] Assay design.

To do this, you can compare your target sequence to databases of sequences and search for regions of sequence similarities. In order to make your assay as specific as possible, regions of similarity can be masked out before submitting your sequence for design, so they are not considered in the assay design. The National Center for Biotechnology Information (NCBI) hosts a database of all published nucleotide and protein sequences. BLAST[®], a sequence comparison algorithm, is available to facilitate nucleotide and protein searching of the NCBI public databases.

A. **How to use BLAST[®] to search for Sequence Similarity**

1. **Submitting your sequence / Starting your query**

- Go to the [NCBI BLAST[®] site](#)
- Choose “Nucleotide-nucleotide BLAST (blastn)” under **Nucleotide**.
- You may choose to BLAST some or all of your cDNA sequence. If you are only interested in a particular region of a transcript, then choose about 300 – 600 bases in that area to BLAST. If you are not sure about where you want the assay located, or you want options, then you may want to BLAST the whole cDNA sequence to find the best exon boundaries with which to work.
- Enter your sequence into the box provided. You may want to search with your masked sequence generated from RepeatMasker. There are three sequence formats that may be entered into this box. (See pg. 8) For more information on this, click on the word [Search](#) next to the box.
- Choose the appropriate [database](#) to search. When searching with a cDNA sequence for a gene expression assay, you would probably want to search at least the ‘est’ (expressed sequence tags) and the ‘nr’ databases for the species you are working with.
- Under ‘Options for advanced blasting’ you can, among other things, [limit your search to a specific organism](#) using the drop down menu, and opt to [filter](#) your query for low complexity sequences (not necessary if searching with output from RepeatMasker).
- Click on ‘BLAST!’ to submit your search.

2. For more information on how to use BLAST®

NCBI has extensive help documentation on the NCBI BLAST® website. This includes [FAQs](#) and [Tutorials](#). Included on the Tutorials page are also an [Introduction to Similarity Searches](#) and a [Glossary of Terms](#).

BLAST® Submission

The screenshot shows the NCBI BLAST submission page. At the top, the NCBI logo is on the left, and the text "nucleotide-nucleotide BLAST" is on the right. Below this, there are tabs for "Nucleotide", "Protein", and "Translations", and a link "Retrieve results for an RID".

The main form area contains a text input field for the sequence. A blue circle highlights the "Search" button. A blue arrow points from the text "Information on format of submission sequence. This sequence is in FASTA format" to the sequence input field. The sequence is in FASTA format:

```
>BC032413 Homo sapiens B lymphoid tyrosine kinase, mRNA
CACCTCTGTCTGCTGCCGGCAGAAAAGCCACAAGCCATGAAAACCTGATTGAGATGAGAAGA
ATTCATCTGGGACTGGCTTTTGCCTTAGGATGGTGTGGAAAGTTGCTCGTTGTTCGCTAGG
AGCCTGCTCCACTGTAAGGGTGTCCGGATCTGAAAGACTATGGTGAAACACCACTGAAGC
ATTGCCAAGGATGGGGCTGGTAAAGTAGCAAAAAGCCGGACAAGGAAAAGCCGATCAAAGA
```

Below the sequence input field are "Set subsequence" fields for "From:" and "To:". A blue arrow points from the text "Information on Databases to search" to the "Choose database" dropdown menu, which is currently set to "nr". Other options in the dropdown include "est", "est_human", "est_mouse", "est_others", "gss", "htgs", "pat", "pdb", "month", "alu_repeats", and "dbsts". There are also "Reset query" and "Reset all" buttons.

Below the database dropdown is a section titled "Options for advanced blasting". A blue arrow points from the text "To limit search to a specific organism..." to the "Limit by entrez query" field, which is currently set to "Homo sapiens [ORGN]".

Below that is a section titled "Options for Filtering for low complexity sequences if query sequence has not been masked". A blue circle highlights the "Choose filter" section, which includes checkboxes for "Low complexity" (checked), "Human repeats", "Mask for lookup table only", and "Mask lower case". Below this are "Expect" (set to 10) and "Word Size" (set to 11) fields. At the bottom is an "Other advanced" text input field.

3. BLAST Results

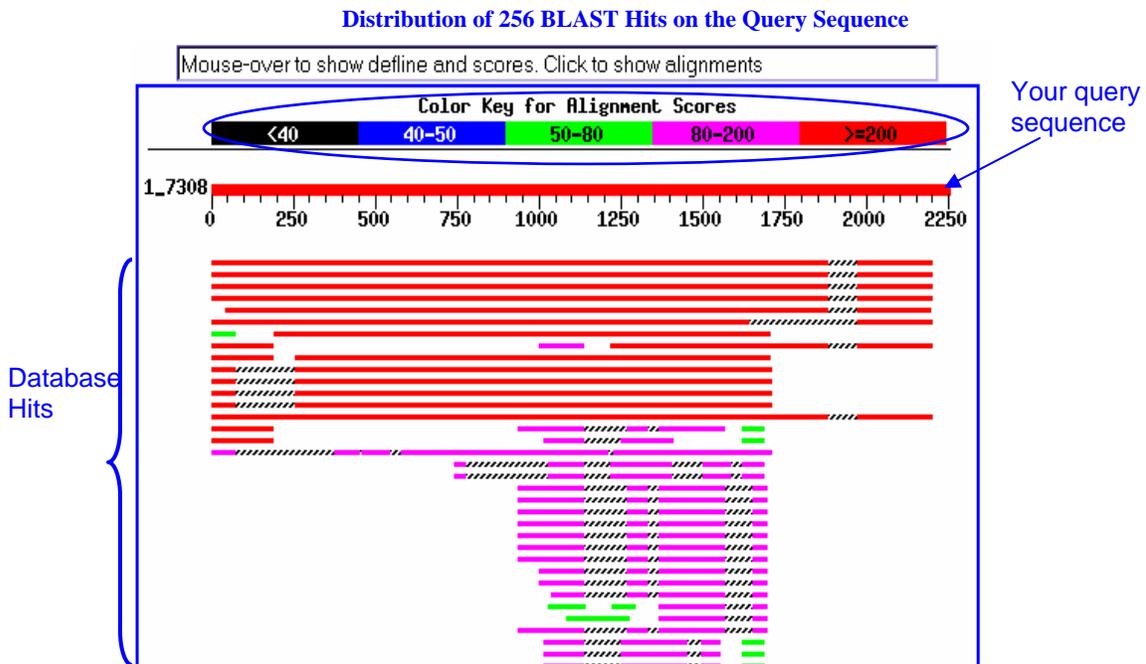
There are three general parts to BLAST® results:

- a. Graphical overview
- b. List of Sequences producing significant alignments to your query
- c. Sequence alignments.

These sections are described below (p9–12) to give you a better understanding of what information can be obtained from a BLAST search of the NCBI public nucleotide database.

a. Graphical Overview

The graphical overview, as seen below, is a representation of the database sequences (hits) that align to your query sequence, with the query sequence represented by the thick red numbered line at the top of the graph. The color of the line represents the score of the alignment, and a striped line connects multiple alignments to the same database sequence.



b. List of Sequences producing significant alignments to your query

The list of sequences is shown from best to worst alignment; the top hit being the best hit (and possibly the sequence with which you queried the database). Public ID information is available as hypertext to the GenBank records that align to your query sequence, as well as a sequence definition. Clicking on the Score hypertext will take you to the actual sequence alignment. The score reflects the degree of similarity between your sequence and the sequence to which it is being aligned. The higher the score is, the more similar the sequences. You should also be able to understand the [E value](#) in order to evaluate the significance of a particular result. The E value represents the number of hits one can "expect" to find by chance when searching a database of a particular size. In this case, the database is the NCBI database that you searched. The lower the E value is, the more significant the match. Hits with E values higher than around 0.1 are unlikely to be very significant.

Click on **Score** to go to sequence alignment

*Sequences producing significant alignments:		Score (bits)	E Value	
gi 21595366 gb BC032413.1	Homo sapiens B lymphoid tyrosine...	3729	0.0	L U
gi 601951 emb Z33998.1 HSBITPTK	H.sapiens mRNA for human ly...	3713	0.0	L U G
gi 33469981 ref NM_001715.2	Homo sapiens B lymphoid tyrosi...	3713	0.0	L U
:				
gi 1015382 gb U34859.1 HSU34859	Human protein tyrosine kina...	373	1e-99	
gi 42557499 gb AF131216.2	Homo sapiens chromosome 8 map 8p...	369	2e-98	
gi 32129354 gb AC090496.28	Mus musculus clone rp23-469n6 m...	186	1e-43	
gi 14140182 emb AJ277921.1 SSC277921	Saimiri sciureus parti...	105	4e-19	
gi 34871626 ref XM_232763.2	Rattus norvegicus Lymphocyte-s...	103	2e-18	L U
gi 33303798 gb AY335586.1	Synthetic construct Homo sapiens...	101	6e-18	
gi 21757951 dbj AK098027.1	Homo sapiens cDNA FLJ40708 fis,...	101	6e-18	L U
:				
gi 5262302 emb AL031729.16 HS159A19	Human DNA sequence from...	62	5e-06	L G
gi 34531137 dbj AK125143.1	Homo sapiens cDNA FLJ43153 fis,...	62	5e-06	U
gi 4885234 ref NM_005248.1	Homo sapiens Gardner-Rasheed fe...	62	5e-06	L U G
gi 182581 gb M12724.1 HUMFGRO7	Human c-fgr proto-oncogene, ...	62	5e-06	L U G
gi 182573 gb M19722.1 HUMFGR	Human fgr proto-oncogene encod...	62	5e-06	L U G

Links to **NCBI databases** for each hit:
LocusLink, UniGene & GEO are shown here

*List shortened for display purposes

By just browsing a list of hits one can get a good idea of the types of sequences that have been found to have some identity to your query. Notice that the first sequence in the list is the one that was used for the search in this example, BC032413.1. The score is very high (3729), and the Expect value is 0. Remember that the closer an E-value is to "0" the more "significant" the match. For this particular query, most of the hits are to human tyrosine kinases, which is the same molecular function as the query. Remember that what you're looking for is the ability to design an assay that will uniquely detect your sequence of interest, whether it is a unique gene sequence or a unique splice variant. If you find some regions of similarity between your sequence and another, those bases can be masked out, so that they will not be considered for assay design.

c. Sequence Alignments

This section is your query sequence aligned to every sequence on your list of hits. These alignments are to help assess the degree of similarity. The Score and Expect values are displayed underneath the sequence identifiers. The number of bases aligned and percent identity are shown, as well as the strand that was aligned of your query sequence and the database hit.

You'll notice that the first hit in this list, shown here boxed in blue, is the query sequences aligned to itself. This will be the first alignment shown, and will be a 100% match to itself.

In the example on page 11, the query sequence did not align to this database hit contiguously. Only a part of the query sequence aligned to this sequence from the database: starting at base 1367 of the query sequence, and ending at base 1565, and from 936 -1136 of the query sequence. Also, there is more than one alignment associated with this hit. These alignments are shown in order of significance, and have different E-values and scores. Sometimes BLAST alignments of an mRNA query can bring up hits to genomic DNA. These alignments can be broken into multiple [HSP segments](#), indicating the presence of introns in the gDNA. If a segment of your query sequence came up with a significant match to part of a sequence from another gene, you could either mask out that region of the sequence in your sequence for submission or simply not include that region in your submission and find another region of interest in your gene to submit.

B. How to use BLAST® dbSNP to search for Sequence Polymorphisms

1. Submitting your sequence / Starting your query

- Go to the [NCBI BLAST® SNP site](#). The default Program is blastn. This is the program you should use.
- Choose the sequence that you would like to submit based on your BLAST for sequence similarity.
- Enter your sequence into the box provided. The sequence format should be [FASTA](#). You may either search with your masked sequence (output from RepeatMasker) or have the sequence filtered for you by the program. To have the sequence filtered for you, simply check the appropriate boxes next to the word [FILTER](#), as shown on page 13.
- Click on 'Submit Query' to submit your search.

Single Nucleotide Polymorphism

Select the BLAST program to use and enter your sequence in the text area below.

Program:

Query Sequence

Enter your sequence as:

```
>BC032413 Homo sapiens B lymphoid tyrosine kinase
CACCTCTGTCTGCTGCCGGCAGAAAGCCACAAGCCATGAAAAGTATTGA
GATGAGAAGAATTCATCTGGGACTGGCTTTTGGCTTTAGGATGGTGTGGGA
AGTTGCTCGTTGTCGCTAGGAGCCTGCTCCACTGTAAGGGTGTCCGGGATC
TGAAGAGCTATGGTGAACACCCTGAAAGCATTGCCAAGGATGGGGCTGG
TAAAGTAGCAAAAAGCCGGACAAAGGAAAAGCCGATCAAAGAGAAGGACAAAG
GGCCAATGGAGCCCCCTGAAGGTCAGCGCCCAAGACAAAGGACGCCCCGCC
ACTGCCGCCCTGGTTGTCTTCAACCACCTTACTCCTCCACCGCCCGATG
```

Snp Blast Databases(Human)

Chr. 1 Chr. 7 Chr. 13 Chr. 19
 Chr. 2 Chr. 8 Chr. 14 Chr. 20
 Chr. 3 Chr. 9 Chr. 15 Chr. 21
 Chr. 4 Chr. 10 Chr. 16 Chr. 22
 Chr. 5 Chr. 11 Chr. 17 Chr. X
 Chr. 6 Chr. 12 Chr. 18 Chr. Y
 MultiChr. NotOnChr. **All of the Above**

BLAST Search Options

Expect **Descriptions** **Alignments**

0.01 100 100

Filter Low complexity Human repeats Mask for lookup table only

Other advanced options:

2. dbSNP BLAST® Results

The output is typical of BLAST® results, a list of sequences producing significant alignments to your query and the sequence alignments. Notice the Scores and Expect values, as well as the public identifiers. These are all discussed in the section entitled [“List of Sequences producing significant alignments to your query”](#).

Sequences producing significant alignments:	Score (bits)	E Value
gnl dbSNP rs922483_allelePos=161totalLen=577	375	e-101
gnl dbSNP rs2250788_allelePos=500totalLen=998	373	e-100
gnl dbSNP rs2245250_allelePos=201totalLen=401	254	1e-64
gnl dbSNP rs2245232_allelePos=201totalLen=401	228	7e-57
gnl dbSNP rs12386974_allelePos=201totalLen=401	194	1e-46
gnl dbSNP rs6994605_allelePos=201totalLen=401	172	4e-40
gnl dbSNP rs2244938_allelePos=201totalLen=401	105	7e-20
gnl dbSNP rs2244931_allelePos=201totalLen=401	88	2e-14
gnl dbSNP rs13272061_allelePos=201totalLen=401	82	1e-12
gnl dbSNP rs11780851_allelePos=284totalLen=484	52	0.001

Sequence Alignments

By looking for mismatches in the alignment (no hash marks) you will be able to identify documented SNPs. These SNPs should also be masked out (changed to N) in your submission sequence so that no primer or probe is designed over this base.

```
>gnl|dbSNP|rs2250788_allelePos=500totalLen=998
      Length = 998
```

```
Score = 373 bits (188), Expect = e-100
Identities = 189/190 (99%)
Strand = Plus / Plus
```

Your sequence of interest

```
Query: 1 cacctctgtctgctgccggcagaaagccacaagccatgaaaactgattgagatgagaaga 60
      |||
Sbjct: 356 cacctctgtctgctgccggcagaaagccacaagccatgaaaactgattgagatgagaaga 415
```

```
Query: 61 attcatctgggactggccttttgccttaggatgggtgttggagttgctcgttgctcgttagg 120
      |||
Sbjct: 416 attcatctgggactggccttttgccttaggatgggtgttggagttgctcgttgctcgttagg 475
```

```
Query: 121 agcctgctccactgtaaggggtgctcgggatctgaagagctatggtgaaacaccactgaagc 180
      |||
Sbjct: 476 agcctgctccactgtaaggggtgctcgggatctgaagagctatggtgaaacaccactgaagc 535
```

```
Query: 181 attgccaagg 190
      |||
Sbjct: 536 attgccaagg 545
```

Documented SNP in dbSNP.
It is important to mask this base before submission.

III. Identifying Exon Junctions

If you are going to order a gene expression assay, it is important to know where the exon junctions are in the cDNA sequence you are submitting for a Custom TaqMan[®] Assay. The TaqMan[®] MGB probe, when possible, should be designed across an exon-exon boundary in order to exclude the detection of genomic DNA. The exon boundaries should serve as your coordinate(s) in your submission file. The more coordinates you provide, the better your chances of having an assay designed. While you may provide as few as one coordinate, or as many as you'd like, only one assay per sequence will be designed. If you are working with a gene sequence that is in a public database there are many places you may go to find exon information on the web. A few are listed and described below.

A. Entrez at NCBI (National Center for Biotechnology Information)

Entrez is a tool used to query different databases at NCBI. GenBank is a public database of nucleotide sequences (as well as other sequences), that is updated daily. A good number of sequences are annotated with mRNA sequences, so you may be able to find some exon information on your sequence of interest here.

To do this:

- Search the nucleotide database using [Entrez](#) at NCBI.
- Find your gene of interest with a complete coding sequence.
- Click on the GenBank Accession number.

NCBI Entrez Nucleotide Search for human IL10. Results show GenBank Accession U16720 for Human interleukin 10 (IL10) gene, complete cds.

- Scroll down to find the **mRNA** annotated under the **Features** section. This will list the exon start and stop bases for this gene sequence. For example, in this record the first exon starts at base 4057 and ends at base 4221, the second exon is from bases 5088 – 5147, and so on. If you join these exons together, they make up the mRNA sequence. The bases listed correspond to the bases from the DNA sequence shown in that particular GenBank record.

```

1: U16720. Human interleukin...[gi:1041812]

LOCUS       HSU16720                8868 bp    DNA     linear   PRI 28-OCT-1995
DEFINITION  Human interleukin 10 (IL10) gene, complete cds.
ACCESSION   U16720
  
```

```

FEATURES             Location/Qualifiers
     source            1..8868
                     /organism="Homo sapiens"
                     /mol_type="genomic DNA"
                     /db_xref="taxon:9606"
                     /chromosome="1"
     repeat_region     1144..1447
                     /rpt_family="Alu"
                     /rpt_type=dispersed
     mRNA              join(<4057..4221,5088..5147,5438..5590,6601..6666,
                     7742..8868)
     gene              join(4057..4221,5088..5147,5438..5590,6601..6666,
                     7742..7834)
                     /gene="IL10"
     CDS               join(4057..4221,5088..5147,5438..5590,6601..6666,
                     7742..7834)
                     /gene="IL10"
                     /codon_start=1
  
```

B. Ensembl / Vega

The [Vertebrate Genome Annotation](#) (VEGA) is part of the [Ensembl](#) project. The VEGA database is a collection of manually curated genome sequences. The current genomes available for searching are human, mouse and zebrafish.

To search this database for exon information:

- Choose the species of interest.
- Enter in a gene identifier, such as gene name, gene ID, accession number, etc, and click on “Lookup”

The screenshot shows the Vega Human Annotation Browser search page. At the top, there are logos for Vega Human, The Wellcome Trust Sanger Institute, and HAVANA. The main search area is titled "Search Vega" and contains a dropdown menu set to "Gene" and a text input field containing "GPR116". A red "Lookup" button is to the right. Below the search area, there are two sections: "Select a Chromosome to Browse" showing a karyotype with chromosomes 6, 7, 9, 10, 13, 14, 20, and 22, and "Data Entry Points" which lists links for MapView, ContigView, BlastView, and TextView. Below these are example data points for Gene (GNRH2), Transcript (C20orf011), Exon (HLA-C-001), Peptide (C20orf18-011), and Contig (AL035460.15.1.135005).

- Once you get the results, click on the link to the gene of interest.

The screenshot shows the Vega Human TextView search results page. At the top, there are logos for Vega Human and TextView. The navigation bar includes links for Home, Human, BLAST, Export Data, Search, and Feedback. The search area shows "Homo sapiens" selected for species, "Gene" for the search type, and "gpr116" in the search field. The results are displayed in "standard" format. A red "Lookup" button and a "Help" button are visible. Below the search area, there is a message: "1 documents match your query (Documents searched: 12524)". A blue arrow points to the first result, which is "1. Vega Gene: OTTHUMG00000014793". The result description includes: "Vega gene OTTHUMG00000014793 has 3 transcripts: OTTHUMT00000040806, OTTHUMT00000040808, OTTHUMT00000040807". It also lists external identifiers: HUGO: dJ365O12.1, 19030, GPR116; LocusLink: dJ365O12.1, 221395; RefSeq: dJ365O12.1, NM_015234; SWISSPROT: dJ365O12.1, Q8IZF2; Vega_gene: OTTHUMG00000014793, dJ365O12.1. A URL is provided at the bottom: http://vega.sanger.ac.uk:80/Homo_sapiens/geneview?gene=OTTHUMG00000014793&db=core.

- Click on “Transcript information” to view the cDNA sequence, with the exons shown contiguously in alternating blue then black text.

Vega Human GeneView The Wellcome Sanger Institute

Home Human BLAST Export Data Search Feedback

Find [e.g. Gene: GNRH2]

Curated Locus Report

Curated Locus	GPR116 (HUGO ID)
Locus ID	OTTHUMG0000014793 [View in Ensembl]
Version	1 (30 Jul 2002)
Classification	Novel_CDS [Definition]
Genomic Location	View gene in genomic location: 46821735 - 46924076 bp (46.8 Mb) on chromosome 6 This gene is located in sequence: AL096772.5.1.117636
Description	No description
Remarks	No remarks
Author	This locus was annotated by Havana < vega@sanger.ac.uk >
Database Matches	HUGO: GPR116 LocusLink: 221395
	<p>1: dJ365O12.1-001 (OTTHUMT00000040806) [Transcript information] [Exon information] [Protein information]</p> <p>2: dJ365O12.1-002 (OTTHUMT00000040807) [Transcript information] [Exon information] [No translation]</p> <p>3: dJ365O12.1-003 (OTTHUMT00000040808) [Transcript information] [Exon information] [No translation]</p>

Vega Human TransView

Transcript cDNA Sequence

```

1 AGAGAACAGAAAGGCAGTTCACCTCTGCTCCCGACAGCCTGGGAAACCCGCAAGAGCCCCAG
61 CATTGAAAGTCTGGTCTTGTGAAACCCACCCCTCCTCTGGCTGTGTGATTGAAATGGGATG
121 CCCTCGAGGTTCACTCACCCTGAGAGGGTTTTGGGCAGATCAGCAGTAAAGGTGTTAAATTT
181 TTAGAAAGCCTGAAAACCTCCAGAAAGAGAAAAGATGAAAATCCCCAAGGAGAAACCACTTTGTGC
241 CTCATGTTTATTGTGATTTATCTTCCAAAAGCTGCACCTGAACTGGAATTAACGAGTCTACT
301 ATTCATCCTTTGAGTCTTCATGAACATGAACCCAGCTGGTGAAGAGGCCTGAGGCAGAAAAA
361 CGAGCCGTTGCCACAAAAAGTCCCTACGGGTGAAGAAATACACTGTAAATATTGAGATCAGT
421 TTTGAAAATGCATCCTTCCCTGGATCCCTATCAAAAGCCTACTTGAACAGCCTCAGTTTCCCA
481 ATTCATGGGAAATAAACAATGACCAAAATACCGACATTTTGGAGATAAATGTGACAAACAGT
541 TGCAGACCTGCTGGAAAATGAAAATCTGGTCTCCTGCGAGACAGGTTATGGGTGGCCTCGG
601 GAAAAGTGTCTTCAAAATCTCATTGTCAAGAGCGGTGACGCTCTCCTCCAGGGCAACCAT
661 TGCAAGTTGCCTTAAAGAACTGCCTCCCAATGGACCTTTTTGCCTGCTTCAGGAAAGATGTT
721 ACCCTGAAACATGAGAGTCAAGTAAATGTAGGCTTTCAAGAAAGCCTCATGAAACACTTCC
781 TCCGCCCTCTATAGGTCCTACAAGACCGACTTGGAAACAGCGTTCCGGAAAGGGTTACGGGA
841 ATTTTACCAGGCTTCAAGGGCGTGACTGTGACAGGGTTCAAAGTCTGGAAAGTGTGGTTGTG
901 ACATATGAAGTCAAGACTACCCACCTACCTTGAGTTAATACATAAAGCCAATGAAACAA
961 GTTGACAGAGCCTCAATCAGACCTCAAAAATGGACTACAACCTCCTTCAAGCAGTTACT
1021 ATCAATGAAAGCAATTTCTTTGTCAACCCAGAAAATCATCTTTGAAAGGGGACACAGTCAGT
1081 CTGGTGTGTGAAAAGGAAGTTTTGTCTCCAATGTGCTTTGGCGCTATGAAAGAACAGCAG
1141 TTGGAAAATCCAGAAACAGCAGCAGATTCTCGATTTACACCGCACTTTTCAACAACATGACT
1201 TCGGTGTCCAAGCTCACCATCCACAACATCACTCCAGGTGATGCAAGTGAATATGTTTGC
1261 AAACGTATATTAGCATTTTTGAATATGAGTGAAGAAAGAAAATAGATGTTATGCCCATC
1321 CAAATTTGGCAAAATGAAAGAAATGAAGGTGATGTGCGCAACAATCCTGTATCTTTGAAC
1381 TGCTGCAAGTCAAGGTAATGTTAATTTGGAGCAAAAGTAAAGTGAAGCAGGAAAGGAAAAATA
1441 AATATCCAGGAACCCCTGAGACAGACATAGATTCTAGCTGAGCAGATACACCCCTCAAG
1501 GCTGATGGAACCCAGTGGCCAAAGCGGGTCTGTTGGAACAACAGTCATCTACACTTTGTGAG
1561 TTCAATCAGTCCATGAGGCAAGGAGGAGTGAACAAATGAAAGTACATTCATCTCTGTG
1621 GCAATCTAACAATAAACCOCGGAACCAATTTCTGTTCTGAGGGACAAAACTTTTCTATA
1681 AAAATGATCAGTGTGAGTAACTATGATGAGGTTTATGGAAACACTTCTGTGGAATTT
1741 AAAATATACCAAAAGATTTTATACCAAGGAGGATCTCTGATGGAGCAGAAATCAGTACTG
1801 ACAGTCAAGACCTCGACAGGGAGTGAATGGAACCTATCACTGCATATTTAGATATAAG
1861 AATTCATACAGTATTGCAACCAAGACGTCATTTGTCACCCGCTGCCTCTAAAGCTGAAC
1921 ATCATGTTGATCCTTTGGAAGCTACTGTTTCAATGCAAGTGGTCCCATCAGTCAAGTGC
1981 TGCAATAGAGGAGGATGGAGACTACAAGTTACTTTCCATACGGGTTCTCATCCCTCTC
2041 GCTGCAAAAGAGCTTAACAAAAACAAGTGTGCTACAACCAAAATTTCAATGCAAGCTCA
2101 GTTCTGTTGTTCAAAAACCTGTTGATGTTGTTGTCCTTTAACAATGCTGCTAATAAT
2161 TCAGTCTGAGCCCATCTATGAAGCTGAATCTGGTCTCGGGGAAAACATCAGTGCAG
2221 GATCCCGTAATAGTGTGCGAGAGCCGGGAAAAGTCATCCAGAACTATGCCGTTCTCA

```

exon 1
 exon 2
 exon 3
 exon 4
 :
 :

Note: If there is no exon information available for your sequence of interest, you may still submit that sequence for assay design. For your coordinates, select multiple sites across the sequence to ensure optimal design.

Having evaluated the quality of your sequence information, you are now ready to move on to preparing your submission file using the [File Builder software](#).

For Research Use Only. Not for use in diagnostic procedures.

Custom TaqMan Gene Expression products –

Notice to Purchaser: Disclaimer of License for Custom Sequence Detection Primers

This product is optimized for use in the Polymerase Chain Reaction (PCR) and 5' nuclease detection methods covered by patents owned by Roche Molecular Systems, Inc. and F. Hoffmann-La Roche Ltd. No license under these patents to use the PCR process or 5' nuclease detection methods is conveyed expressly or by implication to the purchaser by the purchase of this product. A license to use the PCR process for certain research and development activities accompanies the purchase of certain Applied Biosystems reagents when used in conjunction with an authorized thermal cycler, or is available from Applied Biosystems. Further information on purchasing licenses to practice the PCR process may be obtained by contacting the Director of Licensing, Applied Biosystems, 850 Lincoln Centre Drive, Foster City, California 94404 or at Roche Molecular Systems, Inc., 1145 Atlantic Avenue, Alameda, California 94501, USA.

Notice to Purchaser: Disclaimer of License for Custom TaqMan Probes

This product is optimized for use in the Polymerase Chain Reaction (PCR) and 5' nuclease detection methods covered by patents owned by Roche Molecular Systems, Inc. and F. Hoffmann-La Roche Ltd. No license under these patents to use the PCR process is conveyed expressly or by implication to the purchaser by the purchase of this product. A license to use the PCR process for certain research and development activities accompanies the purchase of certain Applied Biosystems reagents when used in conjunction with an authorized thermal cycler, or is available from Applied Biosystems. Further information on purchasing licenses to practice the PCR process may be obtained by contacting the Director of Licensing, Applied Biosystems, 850 Lincoln Centre Drive, Foster City, California 94404 or at Roche Molecular Systems, Inc., 1145 Atlantic Avenue, Alameda, California 94501, USA.

Notice to Purchaser

TaqMan® probes are covered by U.S. Patent 5,723,591 and foreign counterparts and patents pending owned by Applera Corporation, and may be covered by U.S. Patents 5,801,155 and 6,084,102 and foreign counterparts licensed to Applied Biosystems.

Applied Biosystems, Assays-by-Design and ABI PRISM are registered trademarks and AB (Design), Applera, myScience are trademarks of Applera Corporation or its subsidiaries in the U.S. and/or certain other countries.

TaqMan is a registered trademark of Roche Molecular Systems, Inc.
BLAST is a registered trademark of the National Library of Medicine.
All other trademarks are the sole property of their respective owners.

Publication 127GU07-02

Part Number 4371002 Rev B