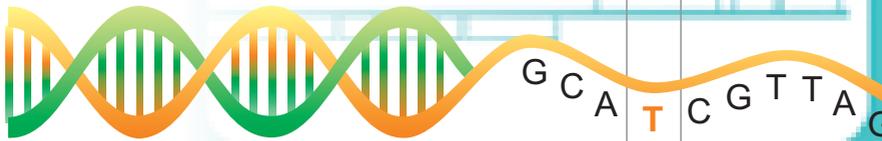
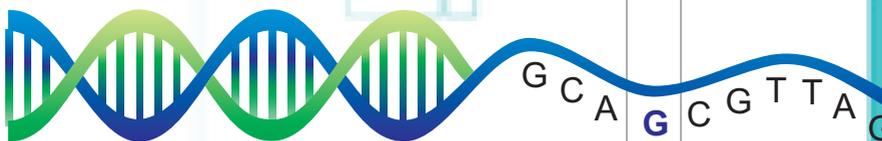


## AnotaSNP: software para organização de informações de anotação de genes em estudos de associação genômica ampla



# ANNOTATION



*Empresa Brasileira de Pesquisa Agropecuária  
Embrapa Informática Agropecuária  
Ministério da Agricultura, Pecuária e Abastecimento*

# **Documentos 129**

## **AnotaSNP: software para organização de informações de anotação de genes em estudos de associação genômica ampla**

*Roberto Hiroshi Higa  
Gabriel Bueno de Oliveira*

## **Embrapa Informática Agropecuária**

Av. André Tosello, 209 - Barão Geraldo  
Caixa Postal 6041 - 13083-886 - Campinas, SP  
Fone: (19) 3211-5700 - Fax: (19) 3211-5754  
www.embrapa.br/informatica-agropecuaria  
sac: www.embrapa.br/fale-conosco/sac/

### **Comitê de Publicações**

Presidente: *Silvia Maria Fonseca Silveira Massruhá*

Secretária: *Carla Cristiane Osawa*

Membros: *Adhemar Zerlotini Neto, Stanley Robson de Medeiros Oliveira, Thiago Teixeira Santos, Maria Goretti Gurgel Praxedes, Adriana Farah Gonzalez, Neide Makiko Furukawa, Carla Cristiane Osawa*

Membros suplentes: *Felipe Rodrigues da Silva, José Ruy Porto de Carvalho, Eduardo Delgado Assad, Fábio César da Silva*

Supervisor editorial: *Stanley Robson de Medeiros Oliveira, Neide Makiko Furukawa*

Revisor de texto: *Adriana Farah Gonzalez*

Normalização bibliográfica: *Maria Goretti Gurgel Praxedes*

Editoração eletrônica/Arte capa: *Neide Makiko Furukawa*

Imagens capa: *disponível em <<https://www.google.com.br/>>*

### **1ª edição**

on-line 2014

#### **Todos os direitos reservados.**

A reprodução não autorizada desta publicação, no todo ou em parte, constitui violação dos direitos autorais (Lei nº 9.610).

#### **Dados Internacionais de Catalogação na Publicação (CIP) Embrapa Informática Agropecuária**

---

Higa, Roberto Hiroshi.

AnotaSNP : software para organização de informações de anotação de genes em estudos de associação genômica ampla / Roberto Hiroshi Higa, Gabriel Bueno Oliveira . - Campinas : Embrapa Informática Agropecuária, 2014.

28 p. : il. ; 14,8 cm x 21 cm. - (Documentos / Embrapa Informática Agropecuária, ISSN 1677-9274; 129.

1. Anotação de snps. 2. Gwas. 3. Estudos de associação genoma amplo. I. Oliveira, Gabriel Bueno. II. Embrapa Informática Agropecuária. Título. II. Série.

CDD 005.15 (21.ed.)

---

© Embrapa 2014

## **Autores**

### **Roberto Hiroshi Higa**

Engenheiro eletricista, doutor em Engenharia Elétrica, Pesquisador da Embrapa Informática Agropecuária, Campinas, SP

### **Gabriel Bueno de Oliveira**

Estudante de Engenharia da Computação, Bolsista CNPq/PIBIC, Campinas, SP

# Apresentação

Este documento apresenta o manual do software anotaSNP, que tem por objetivo apoiar pesquisadores na análise de experimentos envolvendo estudos de associação genótipo-fenótipo. Especificamente, ele foca na fase de análise das anotações dos genes associados ao fenótipo de interesse, a partir do conjunto de SNPs (*Single Nucleotide Polymorphisms*) identificado. O anotaSNP coleta informações sobre as funções biológicas desses genes que, embora públicas, encontram-se dispersas em diferentes sites da internet, apresentando-as de forma organizada para o usuário, na forma de uma página HTML contendo links para os sites de origem da informação, de forma que o pesquisador possa acessá-las de forma ágil e transparente. Espera-se que a automatização do processo descrito e sua inclusão em pipelines de análise de associação em genética animal contribua para uma maior produtividade e eficiência na execução de projetos de pesquisa envolvidos com este tipo de análise, mantendo o foco do pesquisador na interpretação das informações biológicas relacionadas ao estudo.

***Kleber Xavier Sampaio de Souza***

Chefe-geral

Embrapa Informática Agropecuária

# Sumário

<b>Introdução</b> .....	9
<b>Modelo de dados</b> .....	12
<b>Tabelas geneinfo e chromosome</b> .....	12
<b>Tabela taxon</b> .....	13
<b>Tabelas go e gene2go</b> .....	13
<b>Tabela snp</b> .....	13
<b>Tabela snp_panel</b> .....	13
<b>Tabela kegg e gene2kegg</b> .....	14
<b>Tabelas homologo e homolo2gene</b> .....	14
<b>Tabela QTL</b> .....	14
<b>Implementação</b> .....	14
<b>Download e organização dos arquivos: script downloads.py</b> ...	15
<b>Popular banco de dados local: script popula.py</b> .....	15

<b>Atualizar informações start/final position na tabela geneinfo: script position.py</b> .....	16
<b>Consultas</b> .....	16
<b>Obter anotação para a lista de SNPs e genes próximos: script getsnpannotation.py</b> .....	16
<b>Obter anotação para genes em intervalos especificados: script getsinterval.py</b> .....	17
<b>AnotaSNP utilizando a plataforma Galaxy</b> .....	17
<b>Exemplo Ilustrativo</b> .....	18
<b>Conclusão</b> .....	26
<b>Apêndices</b> .....	27
<b>Referências</b> .....	27

# AnotaSNP: software para organização de informações de anotação de genes em estudos de associação genômica ampla

---

*Roberto Hiroshi Higa*  
*Gabriel Bueno de Oliveira*

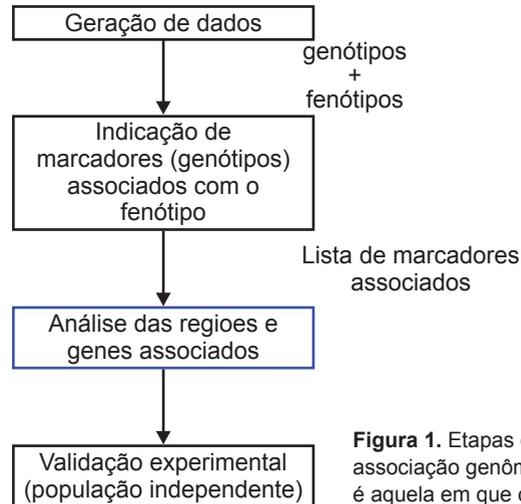
## Introdução

O objetivo dos estudos de análise de associação genômica ampla, *genome-wide association study* (GWAS) é identificar marcadores moleculares do tipo *Single Nucleotide Polymorphisms* (SNP) associados com o fenótipo de interesse (ZIEGLER et al., 2008). Em genética animal, as características analisadas são as de interesse econômico para as cadeias produtivas, como resistência a endo e ectoparasitas, maciez da carne, ausência/ presença de chifres, etc. Esse tipo de análise envolve a utilização de métodos estatísticos e computacionais que manipulam uma grande quantidade de dados para encontrar um conjunto de SNPs que expliquem a variação fenotípica observada.

Em geral, a sequência de etapas envolvidas em um estudo de associação genômica ampla pode ser resumida em:

a) Geração dos dados de genótipo e fenótipo.

- b) Realização da análise de associação genótipo-fenótipo.
- c) Análise das anotações das regiões associadas e identificação de genes candidatos.
- d) Validação experimental utilizando uma população independente (Figura 1).

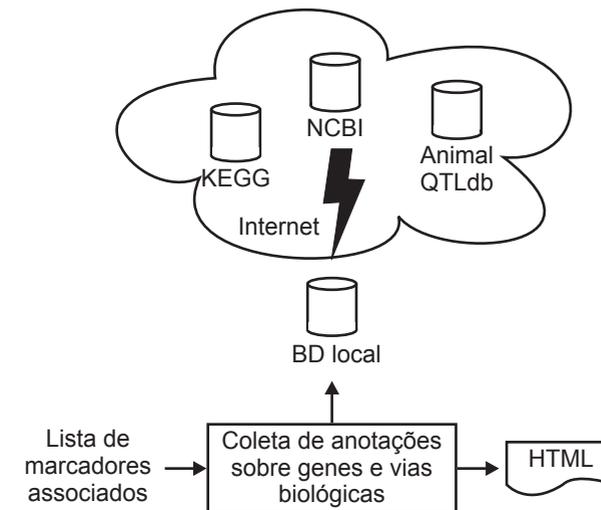


**Figura 1.** Etapas do processo simplificado de análise de associação genômica ampla. A etapa em destaque (azul) é aquela em que o software anotaSNP é utilizado.

A etapa (c) é realizada após a identificação do conjunto de SNPs associados com o fenótipo estudado e tem como primeiro passo mapear esses SNPs no genoma de referência do organismo, identificar os genes presentes em suas vizinhanças e proceder com a coleta de anotações funcionais, como as funções que esses genes desempenham e/ou as vias biológicas de que participam, em diferentes bancos de dados públicos. Normalmente, esse passo é executado de forma manual, acessando diferentes sites da internet, sendo o *Kyoto Encyclopedia of Genes and Genomes* (KEGG) (KEGG, 2013), o National Center for Biotechnology Information (NCBI) (NCBI, 2013), o *Gene Ontology* (GO) (GO, 2013) e o animal QTLdb (ANIMAL QTLdb, 2013), os mais utilizados em estudos sobre genética animal.

Para facilitar esse processo, o software anotaSNP automatiza o processo de coleta de anotações funcionais, realizando o mapeamento dos SNPs no genoma de referência, identificação de genes em sua vizinhança e busca

por informações sobre funções e vias biológicas relacionadas. Para isso, é criada uma base de dados local, considerando informações oriundas das seguintes bases de dados públicas: KEGG, NCBI, GO e Animal QTLdb. O processo automatizado pelo software anotaSNP é composto de duas etapas (Figura 2): a primeira, offline, consiste em fazer o download dos dados de interesse disponíveis nos bancos de dados públicos, KEGG, NCBI e Animal QTLdb, e extrair as informações relevantes, integrando e armazenando-as num banco de dados local. Este banco de dados contém informações sobre SNPs, genes, vias biológicas e anotação funcional, ficando disponíveis para consulta. A segunda etapa, online, consiste em realizar uma consulta tendo como query uma lista de SNPs ou de regiões do genoma de interesse. O resultado é retornado como um arquivo no formato HTML, contendo uma tabela que relaciona os SNPs aos genes próximos, suas respectivas anotações, bem como links para os respectivos sites originais.



**Figura 2.** Processo de coleta de anotações sobre genes e vias biológicas.

Espera-se que a automatização do processo descrito acima e sua inclusão em pipelines de GWAS em genética animal contribuam para uma maior produtividade e eficiência na execução de projetos de pesquisa que demandem este tipo de análise. Dessa maneira, o pesquisador pode manter o foco na interpretação das informações biológicas relacionadas ao estudo.

## Modelo de dados

O banco de dados local implementado pelo software anotaSNP integra as informações oriundas das diferentes fontes de dados públicas (KEGG, NCBI e Animal QTLdb), visando facilitar a construção de resultados de consulta (página HTML) às anotações funcionais de SNPs, genes e vias biológicas associadas com o fenótipo estudado. O correspondente modelo de dados é apresentado na Figura 3.vv

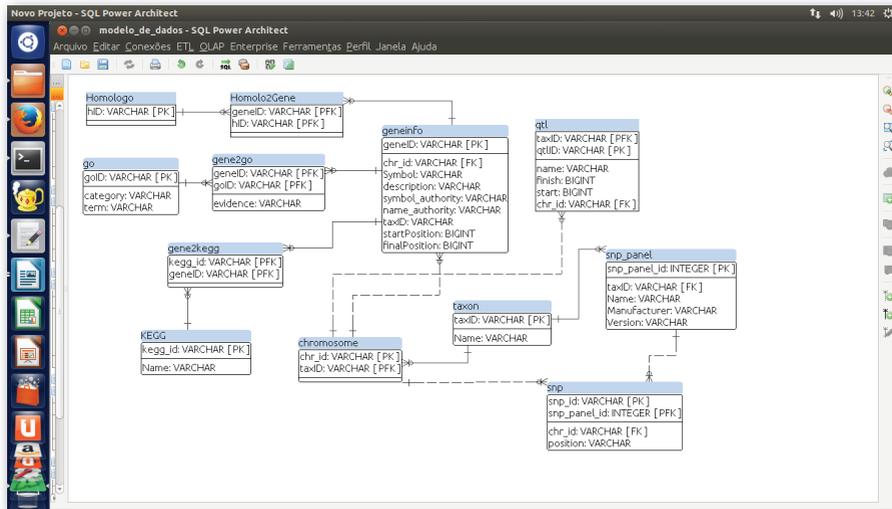


Figura 3. Modelo de dados.

## Tabelas geneinfo e chromosome

A tabela geneinfo contém informações relativas aos genes como: identificação do gene, símbolo associado, descrição, nome, táxon do organismo de origem (táxon é uma unidade taxonômica associada a um sistema de classificação taxonômica, podendo indicar níveis de um sistema de classificação, por exemplo reino, gênero, espécie, subespécie, organismo, etc.), cromossomo onde o gene se encontra, bem como as posições de início e fim. A tabela chromosome contém as informações de identificação do cromossomo e táxon do organismo ao qual ele pertence. Todas as informações são obtidas do arquivo geneinfo, trazido do sítio do NCBI.

## Tabela taxon

A tabela táxon contém o táxon ID e o nome correspondente aos organismos com informações presentes no banco de dados local (são incluídos todos os organismos presentes no arquivo homologene.data).

## Tabelas go e gene2go

As tabelas go e gene2go contêm informações referentes às anotações de GO, como identificação, evidência, categoria e termo. Posteriormente, essas informações são relacionadas aos genes. Essas informações são obtidas do arquivo gene2go, trazido do sítio do NCBI.

## Tabela snp

A tabela snp contém as informações relacionadas aos SNPs, como identificação, cromossomo onde se encontra, posição no cromossomo e organismo a que se refere. Atualmente, a tabela contém dados dos seguintes organismos: *Bos taurus*, *Sus scrofa*, *Ovis aries* e *Gallus gallus*. As informações contidas nessa tabela são obtidas de arquivos específicos, disponibilizados por fabricantes de painéis comerciais, e estão disponíveis no diretório db (diretório incluído no pacote de distribuição do software anotaSNP). São eles: Bovine3K\_2900\_ED2\_13Jul10\_69831.design.csv, BovineSNP50\_v2\_C.csv, SNP\_Map\_cattle.txt, SNP\_Map\_suine.txt, PorcineSNP60v2.txt e Ovine\_HapMap\_MarkerList\_16Jan09.csv.

## Tabela snp\_panel

A tabela snp\_panel contém informações sobre painéis de SNPs, específicos para cada espécie, utilizados para genotipagem. Cada painel cadastrado nessa tabela agrega um conjunto de SNPs descrito na tabela snp.

## Tabela kegg e gene2kegg

As tabelas kegg e gene2kegg contêm informações sobre as vias metabólicas descritas no banco de dados KEGG, relacionando-as com os genes associados. Essas informações são obtidas dos arquivos contidos na pasta est.kegg.jp, trazidos do sítio do KEGG.

## Tabelas homologo e homolo2gene

As tabelas homologo e homolo2gene contêm informações para relacionar os genes homólogos entre os organismos atualmente considerados (*B. taurus*, *S. scrofa*, *O. aries*, *G. gallus*, *Homo sapiens* e *Mus musculus*). Essas informações são consideradas para agregar anotações de genes de organismos modelo como o *M. musculus* e/ou mais estudados como o *H. sapiens* às anotações específicas do organismo estudado. Esses dados são obtidos do arquivo homologene.data, trazido do sítio do NCBI.

## Tabela QTL

A tabela QTL contém as informações sobre *Quantitative Trait Locus (QTLs)* – regiões do genoma relacionados a genes que influenciam a manifestação de caractere quantitativo, por exemplo qualidade de carne em bovinos), incluindo sua designação (nome), cromossomo em que se encontra e as respectivas posições inicial e final da região. Essas informações são obtidas dos arquivos contidos na pasta QTLdb, trazidas do sítio do Animal QTLdb.

## Implementação

Para implementação do software anotaSNP, foi utilizada a linguagem de programação Python (PYTHON, 2013). O banco de dados local foi implementado utilizando o sistema gerenciador de banco de dados (SGBD) PostgreSQL (POSTGRESQL, 2013), acessado por meio do adaptador para python Psycopg2 (PSYCOPG2, 2013). Para manipular o arquivo de anotação de genomas, utilizou-se a biblioteca específica para manipulação de dados biológicos Biopython (BIOPYTHON, 2013).

O software anotaSNP é constituído pelo conjunto de programas python descritos abaixo, que devem ser executados conforme a ordem de precedência indicada.

## Download e organização dos arquivos: script downloads.py

**Execução:** `$ python download.py`

Este script faz download dos arquivos de anotação de genoma e da base de dados nos sites KEGG, NCBI e Animal QTLdb, organizando-os em subdiretórios do diretório <diretório de instalação do anotaSNP>/db. Caso os arquivos já existam no diretório, o script apaga todos os arquivos, baixando-os novamente. Este script, junto com os scripts popula.py e position.py constituem a parte off-line do software.

Os arquivos do sítio Animal QTLdb precisam ser baixados manualmente. Mensagens na parte final do script alertam o usuário para a execução desta tarefa.

## Popular banco de dados local: script popula.py

**Execução:** `$ python popula.py -b <banco> -u <usuario1> -s <senha_banco> [-h p/ ajuda]`

Este script cria o banco de dados local e utiliza os arquivos trazidos pelo script download.py para o diretório <diretório de instalação do anotaSNP>/db para extrair os dados necessários para populá-lo, conforme ilustrado na Figura 1. Caso o banco de dados local já exista, ele é apagado e populado novamente. Este script, combinado aos scripts download.py e position.py constituem a parte offline do software.

As informações sobre os SNPs que compõem painéis comerciais específicos, descritos na Tabela 1, são distribuídas com o software e podem ser encontradas no diretório <diretório de instalação do anotaSNP>/db.

<sup>1</sup> O usuário PostgreSQL especificado deve possuir permissão de *superuser*.

**Tabela 1.** Arquivos contendo informações específicas sobre painéis comerciais de genotipagem.

ID	Painel	Arquivo
1	Bovine HD	SNP_Map_cattle.txt
2	Bovine 50k	BovineSNP50_v2_C.csv
3	Bovine 3k	Bovine3K_2900_ED2_13Jul10_69831.design.csv
4	Swine 60k v1	SNP_Map_swine.txt
5	Swine 60k v2	PorcineSNP60v2.csv
6	Ovine 60k	Ovine_HapMap_MarkerList_16Jan09.csv

## Atualizar informações start/final position na tabela geneinfo: script position.py

**Execução:** `$ python position.py -b <banco> -u <usuario> -s <senha_banco> [-h p/ ajuda]`

Este script faz o parsing (processo de análise da sequência de uma entrada, segundo uma gramática pré-definida, para identificar a estrutura da informação que ela contém) dos arquivos de anotação de cada genoma, trazidos do sítio do NCBI, e extrai informações da localização de cada gene (*startposition* e *finalposition*) cadastrado na tabela geneinfo. Combinado aos scripts `download.py` e `popula.py`, constituem a parte offline do software.

## Consultas

Os scripts que executam consultas (`getspannotation.py` e `getinterval.py`) constituem a parte off-line do software.

## Obter anotação para a lista de SNPs e genes próximos: script getsnpannotation.py

**Execução:** `$ python getsnpannotation.py -s <arquivo_Snps> -p <painel> -d <delta> [-g <organismo_homologo1>, ..., <organismo_homologoN>] -g ALL ] -b <banco> -u <usuario> -passwd <password> [-host <host do banco>] [-h p/ ajuda]`

**Saída:** Documento HTML `output_getsnpannotation.HTML`. Se a opção `-g` é utilizada, anotações de genes homólogos dos organismos especificados são agregadas ao documento HTML de saída.

**OBS:** o parâmetro `-p <painel>` é especificado pelo valor da coluna ID na tabela 1.

## Obter anotação para genes em intervalos especificados: script getsinterval.py

**Execução:** `$ python getsinterval.py -s <arquivo_Intervalo> -p <painel> [-g <organismo_homologo1>, ..., <organismo_homologoN>] -g ALL ] -b <banco> -u <usuario> -passwd <password> [-host <host do banco>] [-h p/ ajuda]`

**Saída:** Documento HTML `output_getsinterval.HTML`. Se a opção `-g` é utilizada, anotações de genes homólogos dos organismos especificados, são agregadas ao documento HTML de saída.

**OBS:** o parâmetro `-p <painel>` é especificado pelo valor da coluna ID na Tabela 1.

## AnotaSNP utilizando a plataforma Galaxy

Também é possível realizar consultas sobre anotações de SNPs (script `getsnpannotation.py`) ou regiões genômicas (script `getinterval.py`) providas pelo software `anotaSNP` por meio da plataforma Galaxy (GOECKS et al., 2010). Para isso, é preciso:

- Instalar o `anotaSNP` e popular o banco de dados local, conforme exemplificado nas seções anteriores.
- Incluir os scripts de consulta do `anotaSNP` como ferramentas do Galaxy.

O procedimento abaixo ilustra como incluir o script `getsnpannotation.py` como uma ferramenta no Galaxy:

1. Criar um diretório para receber o script (ex: `myTools`) dentro do diretório de instalação do Galaxy (ex: `~/galaxy-dist/tools/myTools`).

- Realizar uma cópia dos arquivos getsnpannotation.py e getsnpannotation.xml para o diretório recém-criado.
- No diretório de instalação do Galaxy (~/.galaxy-dist) encontre o arquivo tool\_conf.xml, modificando-o conforme destacado na Figura 4.

```

1 <?xml version='1.0' encoding='utf-8'?>
2 <toolbox>
3 <section name="myTools" id="mTools">
4 <tool file="myTools/getsnpannotation.xml" />
5 <tool file="myTools/getsInterval.xml" />
6 </section>
7 <section id="getText" name="Get Data">
8 <tool file="data_source/upload.xml" />
9 <tool file="data_source/ucsc_tablebrowser.xml" />
10 <tool file="data_source/ucsc_tablebrowser_test.xml" />
11 <tool file="data_source/ucsc_tablebrowser_archaea.xml" />
12 <tool file="data_source/bc_browser.xml" />
13 <tool file="data_source/eft_sra.xml" />
14 <tool file="data_source/microbial_import.xml" />
15 <tool file="data_source/bionart.xml" />
16 <tool file="data_source/bionart_test.xml" />
17 <tool file="data_source/cbt_rice_mart.xml" />
18 <tool file="data_source/granene_mart.xml" />
19 <tool file="data_source/fly_modencode.xml" />
20 <tool file="data_source/flyncine.xml" />
21 <tool file="data_source/flyncine_test.xml" />
22 <tool file="data_source/modncine.xml" />
23 <tool file="data_source/mousemine.xml" />
24 <tool file="data_source/ratncine.xml" />
25 <tool file="data_source/yeastncine.xml" />
26 <tool file="data_source/metabncine.xml" />
27 <tool file="data_source/worm_modencode.xml" />
28 <tool file="data_source/wormbase.xml" />
29 <tool file="data_source/wormbase_test.xml" />
30 <tool file="data_source/eupathdb.xml" />
31 <tool file="data_source/encode_db.xml" />
32 <tool file="data_source/epigraph_import.xml" />
33 <tool file="data_source/epigraph_import_test.xml" />
34 <tool file="data_source/hbvar.xml" />
35 <tool file="genomespace/genomespace_file_browser_prod.xml" />
36 <tool file="genomespace/genomespace_importer.xml" />

```

Figura 4. Alterações no arquivo “tool\_config.xml”.

Um procedimento similar deve ser executado para incluir o script getinterval.py com uma ferramenta no Galaxy, considerando o arquivo getinterval.xml.

## Exemplo Ilustrativo

- Abra o terminal, vá para o diretório que contém os arquivos do software anotaSNP. Em seguida, execute o seguinte comando: python download.py (sugere-se verificar a ocorrência de erros/problemas entre as mensagens emitidas pelo comando). Ao final do processo de download serão indicados 3 sites para o download manual (Figura 5). Realize os downloads dos arquivos solicitados, armazenando-os no diretório db/QTldb.
- Execute o comando python para executar os scripts “popula.py” e “position.py”, como mostra a Figura 6. Ao final do processo o banco de dados

```

gabriel@llicd002:~/workspace/projetoty
TERMINADO --2013-12-02 10:42:19--
Tempo total: 13s
Baixados: 1 arquivos, 516K em 4,3s (120 KB/s)
--2013-12-02 10:42:19-- http://rest.kegg.jp/list/pathway/mmu
Resolvendo proxy.cnpntia.embrapa.br (proxy.cnpntia.embrapa.br)... 10.129.0.245, 20
0.0.70.254
Conectando-se a proxy.cnpntia.embrapa.br (proxy.cnpntia.embrapa.br)[10.129.0.245]:
3128... conectado.
A requisição Proxy foi enviada, aguardando resposta... 200 OK
Tamanho: não especificada [text/plain]
Salvando em: "rest.kegg.jp/list/pathway/mmu"
[ <=> ] 17.515 --K/s em 0,001s
2013-12-02 10:42:20 (11,0 MB/s) - "rest.kegg.jp/list/pathway/mmu" salvo [17515]
TERMINADO --2013-12-02 10:42:20--
Tempo total: 1,2s
Baixados: 1 arquivos, 17K em 0,001s (11,8 MB/s)
--2013-12-02 10:42:20-- http://rest.kegg.jp/link/mmu/pathway
Resolvendo proxy.cnpntia.embrapa.br (proxy.cnpntia.embrapa.br)... 10.129.0.245, 20
0.0.70.254
Conectando-se a proxy.cnpntia.embrapa.br (proxy.cnpntia.embrapa.br)[10.129.0.245]:
3128... conectado.
A requisição Proxy foi enviada, aguardando resposta... 200 OK
Tamanho: não especificada [text/plain]
Salvando em: "rest.kegg.jp/link/mmu/pathway"
[ <=> ] 582.466 72,7K/s em 8,7s
2013-12-02 10:42:37 (65,2 KB/s) - "rest.kegg.jp/link/mmu/pathway" salvo [582466]
TERMINADO --2013-12-02 10:42:37--
Tempo total: 17s
Baixados: 1 arquivos, 569K em 8,7s (65,2 KB/s)
Site para os downloads:
http://www.animalgenome.org/cgi-bin/QTldb/BT/download?file=gbb8tau_4.6
http://www.animalgenome.org/cgi-bin/QTldb/GC/download?file=gbb6_4.6
http://www.animalgenome.org/cgi-bin/QTldb/SS/download?file=gbb5_10.2
gabriel@llicd002:~/workspace/projetoty$

```

Figura 5. Processo de download dos bancos de dados, com indicação de sites para execução manual.

```

gabriel@llicd002:~/workspace/projetoty$ python popula.py -s 123 -b banco -u gabriel
BANCO: banco
USUARIO: gabriel
SENHA: 123
Populando banco de dados 'banco' com o usuario 'gabriel'
Password:
banco de dados 'banco' criado
Conectado ao banco de dados...
Tabelas criadas...
Tabela 't_taxon' Populada...
Tabela 't_chromosome' Populada...
Tabela 't_go' Populada...
Tabela 't_geneinfo' Populada...
Tabela 't_gene2go' Populada...
Tabela snp_panel populada...
Tabela 't_snp' Populada...
Tabela 't_Homologo' Populada...
Tabela 't_Homolo2gene' Populada...
Tabela 't_kegg' Populada...
Tabela 't_gtl' Populada...
Tabela 't_gene2kegg' Populada...
Tabelas populadas com sucesso.
gabriel@llicd002:~/workspace/projetoty$ python position.py -s 123 -b banco -u gabriel
BANCO: banco
USUARIO: gabriel
SENHA: 123
Conectado ao banco de dados...
genome:cattle atualizado...
genome:pig atualizado...
genome:chicken atualizado...
genome:sheep atualizado...
genome:human atualizado...
genome:house_mouse atualizado...
Informacoes start e final position atualizadas.
gabriel@llicd002:~/workspace/projetoty$

```

Figura 6. Execução dos scripts popula.py e position.py.

já estará carregado. Note que esses scripts (download.py, popula.py e position.py) são executados ao se instalar o software anotaSNP, tal que após serem executados, a base de dados está pronta para a realização de consultas (scripts getsnpannotation.py e getsinterval.py).

- Para realizar a consulta usando uma lista de SNPs, basta executar o script getsnpannotation.py (Figura 7).

```
gabriel@l1cd002:~/workspace/projetopy$ python getsnpannotation.py -s Input_Snps.txt -p 1 -d 500000 -b banco -u gabriel -passwd 123 -g 10090
Arg de Snp's: 'Input_Snps.txt'
Painel: 1
Delta: 500000
Organismo homologo: 10090
Banco: banco
Usuario: gabriel
Senha do BD: 123
Arquivo de saída 'output_getsnpannotation.html' criado
gabriel@l1cd002:~/workspace/projetopy$
```

Figura 7. Exemplo ilustrativo da execução do script getsnpsannotation.py.

A página HTML resultante (Figura 8) apresenta as informações coletadas para lista de SNPs query, disponibilizando os links (Figuras 9 e 10) para os sites originais (NCBI, KEGG, GO e Animal QTLdb), onde pode-se obter mais informações sobre os genes/vias biológicas/QTLs próximos dos SNPs query. A página HTML resultante obedece ao seguinte esquema de cores:

- Dourado: início de um novo SNP.
- Laranja: nível superficial onde se encontram os genes e a região de QTL presente no intervalo.
- Amarelo: nível interno a gene, onde se encontram as informações oriundas de KEGG e GO do gene.
- Vermelho: nível interno a gene, onde se encontra a informação de gene homólogo.
- Verde: nível interno a gene homólogo, onde se encontram as informações do KEGG e GO do gene homólogo.
- Azul: links para direcionamento do usuário para fontes de dados específicas (NCBI/Gene, GO, KEGG, NCBI/snpdb e QTLdb).

The screenshot shows a web interface for SNP query results. At the top, it displays 'ChrID: 14', 'SNP ID: ARS-BFGL-BAC-10591', and 'Position: 1754926'. Below this, there is a section titled 'Região de QTL' containing a list of biological processes such as '-Somatic cell score', '-Milk fat percentage', and '-Calving ease (direct)'. A 'Go' button is located below this list. Further down, there is a section titled 'Gene Homólogo: WDYHV motif containing 1' which is highlighted in green. Below this, there is another section titled 'Gene: uncharacterized LOC101907615' which is highlighted in orange. The page is color-coded according to the legend in Figure 7.

Figura 8. Parte superior ilustrando uma página HTML contendo a relação genes/vias biológicas/QTLs próximos aos SNPs query.

The screenshot shows the dbSNP (NCBI) website interface. The search bar contains 'SNP' and the search results show a single hit for 'submitter SNP ID like ARS-BFGL-BAC-10591'. The results table includes columns for Organism, Submitter Handle, Submitter Local SNP ID, NCBI Assay ID(ss#), and RefSNP Cluster ID(rs#). The specific hit is for 'cow\_9913 BFGL\_BARC\_USDA ARS-BFGL-BAC-10591 ss105236242 rs109076417'. The page also features navigation links for various database sections like GENERAL, HUMAN VARIATION, and SEARCH.

Figura 9. Página do sítio dbSNP (NCBI) alcançada a partir de link na página HTML de saída na Figura 7.

The screenshot displays the NCBI Gene page for LOC101907564. The gene is identified as 'uncharacterized LOC101907564 [ Bos taurus (cattle) ]'. The summary section provides details such as Gene symbol (LOC101907564), Gene description (uncharacterized LOC101907564), Gene type (ncRNA), RefSeq status (MODEL), and Organism (Bos taurus). The genomic context section shows the gene's location on Chromosome 14, with coordinates from 17507124 to 17509001. A diagram below the text illustrates the gene's structure on the chromosome, with exons represented by boxes and introns by lines with arrows. The diagram also shows other genes in the region: TMEM5, FKBP5, and HNRG3.

Figura 10. Página de sítio com informações sobre gene alcançada a partir de link na página de saída na Figura 7.

Para realizar uma consulta usando uma lista de intervalos, especificada pelo cromossomo, início e fim do intervalo, executa-se o comando getsinterval.py (Figura 11).

```
gabriel@licd002:~/workspace/projetopy$ python getsinterval.py -s Input_interval.txt -p 4 -g 10090 -b banco -u gabriel -passwd 123
Arq: 'Input_interval.txt'
Painel: 4
Organismo homologo: 10090
Banco: banco
Usuario: gabriel
Senha do BD: 123
Arquivo de saída 'output_getsinterval.html' criado
gabriel@licd002:~/workspace/projetopy$
```

Figura 11. Exemplo de execução do comando getsinterval.py.

A página HTML resultante (Figura 12) apresenta as informações coletadas para a lista de regiões genômicas, disponibilizando links (Figuras 13 e 14) para os sites originais (NCBI, KEGG, GO e Animal QTLdb), onde pode-se

The screenshot shows a web page with a table of genomic regions. The table has columns for 'Início' (Start) and 'Fim' (End) in kb, and 'ChrID'. The first row shows a region on Chromosome 4 from 1.000 kb to 8575.354 kb. Below the table, there is a list of biological pathways and QTLs. Two specific entries are highlighted with arrows: '-Front leg conformation' (labeled as Figura 12) and '-Autoimmune thyroid disease' (labeled as Figura 13). The page also lists several genes, including 'uncharacterized LOC102166976', 'G protein-coupled receptor 20', 'thyroglobulin', 'PHD finger protein 20-like 1', 'uncharacterized LOC102158201', 'maestro heat-like repeat family member 1', 'uncharacterized LOC102161628', 'maestro heat-like repeat-containing protein family member 1-like', 'KH domain containing, RNA binding, signal transduction associated 3', and '5-oxoprolinase (ATP-hydrolyzing)'.

Figura 12. Página HTML contendo a relação genes/vias biológicas/QTLs próximos às regiões query.

**Pig QTL Database - Mozilla Firefox**

Entrada (1) - xgabriel.bueno... Annotation Facebook Pig QTL Database KEGG PATHWAY: ssc05320

www.animalgenome.org/cgi-bin/QTLdb/SS/qdetails?QTL\_ID=16463

**PigQTLdb** Browse Search View Maps F.A.Q.

**QTL #16463 Description:**

**QTL Trait Information**

QTL Trait: <b>Front leg conformation</b>	Vertebrate Trait: <b>Forelimb conformation trait</b>
QTL Symbol: <b>LSCOREF</b>	Product Trait Ontology: n/a
	Clinical Measurement Ontology: n/a

**QTL Map Information**

Chromosome: 4	QTL Peak Location: 8,27 (cM)
QTL Span: n/a	Upper, "Suggestive": n/a
Upper, "Significant": n/a	Peak: rs80811473
Lower, "Significant": n/a	Lower, "Suggestive": n/a

**QTL Experiment in Brief**

**Animals:** Animals were gilts from a Large White line or a Large White x Landrace cross.

**Breeds** that this QTL is associated with: Large white, Landrace

**Design:** Animals were genotyped for 64,232 SNPs using the Porcine 60K BeadChip and analyzed for backfat, loin muscle area, body conformation, and feet and leg structural soundness traits. In total, 51,385 SNPs were used in the analysis.

**Analysis:** SNPs were jointly fitted using Bayesian techniques as random effects in a mixture model that assumed a known large proportion of SNPs had zero effect.

**Software:** GenSel, Haploview

**Notes:**

**Links:** Edit

**Reference**

**Authors:** Fan B, Onteru SK, Du ZQ, Garrick DJ, Stalder KJ, Rothschild MF

Figura 13. Página de site com informações sobre QTLs alcançada a partir de link na página de saída na Figura 11.

**KEGG PATHWAY: ssc05320 - Mozilla Firefox**

Entrada (1) - xgabriel.bueno... Annotation Facebook KEGG PATHWAY: ssc05320

www.genome.jp/dbget-bin/www\_bget?ssc05320

**KEGG PATHWAY: ssc05320** Help

<b>Entry</b>	ssc05320 Pathway
<b>Name</b>	Autoimmune thyroid disease - Sus scrofa (pig)
<b>Description</b>	The classification of autoimmune thyroid disease (AITD) includes Hashimoto's thyroiditis (HT) or chronic autoimmune thyroiditis and its variants, Graves' disease (GD) and autoimmune atrophic thyroiditis or primary myxedema. HT is characterized by the presence of goitre, thyroid autoantibodies against thyroid peroxidase (TPO) and thyroglobulin (Tg) in serum and varying degrees of thyroid dysfunction. During HT, self-reactive CD4+ T lymphocytes (Th) recruit B cells and CD8+ T cells (CTL) into the thyroid. Disease progression leads to the death of thyroid cells and hypothyroidism. Both autoantibodies and thyroid-specific cytotoxic T lymphocytes (CTLs) have been proposed to be responsible for autoimmune thyrocyte depletion. In GD, the TSH-R is the most important autoantigen. Antibodies directed against it mimic the effects of the hormone on thyroid cells, TSH, stimulating autonomous production of thyroxine and triiodothyronine and causing hyperthyroidism. The presence of TSH-R-blocking antibodies that bind the TSH receptor in a similar fashion to the antibodies in patients with grave's disease but that block rather than activate the receptor explains some cases of atrophic hypothyroidism.
<b>Class</b>	Human Diseases; Immune diseases <a href="#">BRITe Hierarchy</a>
<b>Pathway map</b>	ssc05320 Autoimmune thyroid disease

**ALL Links**

- Pathway (1)
- BioSystems (1)
- Genome (1)
- KEGG GENOME (1)
- Gene (53)
- KEGG GENES (53)
- All databases (55)

Download RDF

Figura 14. Página de sítio com informações sobre vias biológicas do KEGG, alcançada a partir de link na página de saída na Figura 11.

obter mais informações sobre os genes/vias biológicas/QTLs próximos das regiões query. A página obedece ao seguinte esquema de cores:

- Dourado: início de um novo intervalo.
- Laranja: nível superficial onde se encontra os genes e a região de QTL presente no intervalo.
- Amarelo: nível interno a gene, onde se encontram as informações do KEGG e GO do gene.
- Vermelho: nível interno a gene onde se encontra a informação de gene homólogo.
- Verde: nível interno a gene homólogo, onde se encontram as informações do KEGG e GO do gene homólogo.
- Azul: Links para direcionamento do usuário para fontes de dados específicas (NCBI/Gene, GO, KEGG, NCBI/snpdb e QTLdb).

Para obter anotações referentes a uma lista de SNPs por meio da plataforma Galaxy, e pressupondo que os procedimentos indicados na seção "AnotaSNP utilizando a plataforma Galaxy" foram executadas com sucesso, basta executar o Galaxy normalmente, abrir a ferramenta referente ao anotaSNP, inserir o arquivo de entrada, as informações solicitadas, e executar o procedimento (Figura 15).

**Galaxy** Analyze Data Workflow Shared Data Visualization Help User Using 4.6 KB

Tools containing all relevant information for each SNP

Annotation Interval Create HTML containing all relevant information for each interval

**Get Data**

- Upload File from your computer
- UCSC Main table browser
- UCSC Test table browser
- UCSC Archive table browser
- BX table browser
- EBI SRA ENA SRA
- Get Microbial Data
- BioMart Central server
- BioMart Test server
- CBI Rice Marj rice mart
- Gramene/Marj Central server
- modENCODE fly server
- Flymine server
- Flymine test server
- modENCODE mod/mine server
- MouseMine server
- Ratmine server

**Annotation (version 1.0.0)**

**File in:** Entrada.txt

**Panel number( 1:bovineHD, 2:bovine50K, 3:bovine3K, 4:Suine60k, v1, 5:Suine60k, v2, 6:Ovine50K):** 1

**Database name:** banco

**User name:** xgabriel

**Password:** xgabriel

**Delta (bp):** 25000

**Database server:** ssc05320

**Taxon of homologous organisms (ex: tx1, tx2, ...):** ALL

Execute

History Unnamed history 4.6 KB

1: Entrada.txt

Figura 15. Exemplo de execução do anotaSNP utilizando a plataforma Galaxy.

Após a execução, basta clicar o ícone *view data* para ter acesso às informações de anotação em formato HTML (Figura 16). Um procedimento semelhante deve ser seguido para obter anotações referentes a uma lista de regiões genômicas.

Figura 16. Exemplo de saída anotaSNP, acessado por meio da plataforma Galaxy.

## Conclusão

O software anotaSNP apoia análises de associação genótipo-fenótipo em espécies animais de interesse agropecuário (bovinos, suínos e aves), reunindo em uma página HTML informações publicamente disponíveis sobre genes próximos ao conjunto de genes obtidos da análise de associação. Este manual apresenta o modelo de dados e o conjunto de scripts que o compõem, bem como exemplos ilustrativos que permitem ao usuário (pesquisadores) utilizá-lo para obter informações biológicas associadas ao conjunto de SNPs obtido da análise de associação genótipo-fenótipo

## Apêndices

- SNPs utilizados no exemplo apresentado:

ARS-BFGL-BAC-10591  
ARS-BFGL-BAC-12452  
ARS-BFGL-BAC-12987  
UA-IFASA-9685

- Intervalo utilizado nos exemplos (Cromossomo!Início do intervalo!Final do intervalo):

4!1000!8575354  
5!4500000!9000000  
14!3000000!6000000

## Referências

- ANIMAL QTLdb. 2013. Disponível em: <<http://www.animalgenome.org/cgi-bin/QTLdb/index>>. Acesso em: 27 set. 2013.
- BIOPYTHON. **Biopython library**. Disponível em: <[Biopython.org](http://Biopython.org)>. Acesso em: 23 set. 2013.
- GOECKS, J.; NEKRUTENKO, A.; TAYLOR, J. Galaxy: a comprehensive approach for supporting accessible, reproducible and transparent computational research in the life sciences. **Genome Biology**, v. 11, p. R86, Aug. 2010. Doi:10.1186/gb-2010-11-8-r86.
- KEGG. **Kyoto Encyclopedia of Genes and Genomes**. 2013. Disponível em: <<http://www.genome.jp/kegg/>>. Acesso em: 27 set. 2013.
- NCBI. **NCBI Resources Gene database**. 2013. Disponível em: <<http://www.ncbi.nlm.nih.gov/gene>>. Acesso em: 27 set. 2013.
- POSTGRESQL. **PostgreSQL**: the world's most advanced open source database. Disponível em: <<http://www.postgresql.org/>>. Acesso em: 27 set. 2013.
- PSYCOPG2. **Psycopg2 tutorial**. Disponível em: <[http://wiki.postgresql.org/wiki/Psycopg2\\_Tutorial](http://wiki.postgresql.org/wiki/Psycopg2_Tutorial)>. Acesso em: 27 set. 2013.
- PYTHON. **Python programming language – official website**. Disponível em: <<http://www.python.org/>>. Acesso em: 27 set. 2013.

THE GENE ONTOLOGY CONSORTIUM. 2013. Disponível em: <<http://geneontology.org/>>. Acesso em: 27 set. 2013.

ZIEGLER, A.; KÖNIG, I.R.; THOMPSON, J. R. Biostatistical aspects of genome-wide association studies. **Biometrical Journal**, v. 50, n. 1 p. 8-28. 2008. DOI: 10.1002/bimj.200710398.



---

*Informática Agropecuária*



Ministério da  
Agricultura, Pecuária  
e Abastecimento

