# Find Your Own Bioinformatics

**Representor :  Xu Wang**

**Group Members :   Liyu  Huang**
**Xu  Wang**
**Yilong  Yang**

**Institute of Crop Sciences**
**Chinese Academy of Agricultural Sciences**

# What's for my own 'Bioinformatics' ??

❧ Text mining- finding the nuggets in the literature

- ◆ iHOP
- ◆ GOPubMed

❧ Before my starting to clone genes …

- ◆ GENEVESTIGATOR
- ◆ Diurnal
- ◆ Codontree

…

❧ Bioinformatics in Plant Biology

*From : Annu. Rev. Plant Biol. 2006. 57:335–60*

**❧ Text mining- finding the nuggets in the literature**

*Q : What can text-mining offer us ??*

✓ 'The goal of text mining is to allow researchers to identify needed information and shift the burden of searching from researchers to the computer. ' *(Rhee, 2006)*

*Q : What do I use text-mining tools for ??*

✓ searching my interested paper with more ease

✓ extracting the useful data from PubMed

✓ ...

**ॐ Text mining- finding the nuggets in the literature**

**Tools I often used :**

◆ iHOP- a gene network for navigating the literature

✓ iHOP provieds the network as a natural way of accessing millions of PubMed abstracts. By using genes and proteins as hyperlinks between sentences and abstracts, the information in PubMed can be converted into one navigable resource, bring all advantages of the Internet to scientiifc literature research.
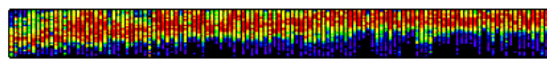
http://www.ihop-net.org/UniPub/iHOP/

| Symbol | Name | | Synonym/ DB-reference | Organism | Results |
|--------|------|--|----------------------|----------|---------|
| | | Life cycles of successful genes | | | |
| SUM1 | SUM1 (SMALL UBIQUITIN-LIKE MODIFIER 1) | | SUMO1 | Arabidopsis thaliana | |

**Search Gene**

**Show overview** new
**Find in this Page**

| Symbol | Name | Synonyms | Organism |
|--------|------|----------|----------|
| SUM1 | SUM1 (SMALL UBIQUITIN-LIKE MODIFIER 1) | At4g26840, F10M23_180, F10M23.180, SMALL UBIQUITIN-LIKE MODIFIER 1, SMT3, SUMO1, SUMO 1, Ubiquitin-like protein SMT3 | Arabidopsis thaliana |

| | |
|--|--|
| UniProt | P55852, Q547B9, Q9SZ24 |
| NCBI Gene | 828791 |
| NCBI RefSeq | NP_194414 |
| NCBI RefSeq | NM_118818 |
| NCBI UniGene | 828791 |
| NCBI Accession | AAP37796, AAL62360 |

more than **1,500 organisms. 80,000 genes. 12 million sentences.**
**...always up-to-date.**

**Homologues of SUM1 ...**

**Definitions for SUM1** ...

**Most recent information for SUM1** ... new
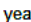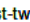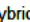
**Enhanced PubMed/Google query ...**

WARNING: Please keep in mind that gene detection is done automatically and can exhibit a certain error. Read more about synonym ambiguity and the iHOP confidence value.

**Find in this Page**

**Sentences in this view contain interactions of SUM1 - Interaction Information is available whenever you see this symbol** - **Read more.**
**For a summary overview of the information in this page click here.** new

Show all
Order by relevance

In yeast-two hybrid assays, **AtSUMO1**/2 **interacts** specifically with a SUMO-conjugating enzyme but not with a **ubiquitin-conjugating enzyme**. [2003]

The levels of **SUMO1** and -2 conjugates but not **SUMO3** conjugates increased substantially following exposure of seedlings to stress conditions, including heat shock, **H(2)O(2)**, **ethanol**, and the amino acid analog **canavanine**. [2003]

**Small ubiquitin-like modifier** (**SUMO**) is a small protein that is structurally related to but functionally different from **ubiquitin [?]**. [2003]

**Small ubiquitin-like modifier** (**SUMO**) is a member of the superfamily of **ubiquitin**-like polypeptides that become covalently attached to various intracellular target proteins as a way to alter their function, location, and/or **half-life**. [2003]

We report the identification and functional analysis of **AtSUMO1**, **AtSUMO2**, and AtSCE1a as components of the **SUMO** conjugation (sumoylation) pathway in **Arabidopsis**. [2003]

Analysis of **transgenic plants** showed that overexpression of **AtSUMO1**/2 does not have any obvious effect in general plant development, but increased sumoylation levels attenuate **abscisic acid** (ABA)-mediated growth inhibition and amplify the induction of **ABA**- and stress-responsive genes such as **RD29A**. [2003]

**ESD4** shows a similar **function** to these proteases **in vitro** and processes the precursor of **Arabidopsis** SUMO (**AtSUMO**) to generate the mature form. [2003]

This activity of **ESD4** is prevented by mutations that **affect** the predicted **active site** of the protease or the cleavage site of the **AtSUMO** precursor. [2003]

This is suggested because **esd4** mutants contain less free **AtSUMO** and more **SUMO** conjugates than wild-type plants, and a **transgene** expressing mature **SUMO** at high levels **enhanced** aspects of the **esd4 phenotype**. [2003]

**ESD4** defines an important role for protein modification by **AtSUMO** in the regulation of flowering. [2003]

**Text mining- finding the nuggets in the literature**

Tools I often used :

◆ GOPubMed- Gene Ontology and PubMed

http://www.gopubmed.org/

∽ **Before my starting to clone genes …**

**Some helpful considerations :**

✓ Is your gene expressed with tissue-specificity ??

✓ Is your gene expressed throughout the life cycle ??

✓ How to determine the sampling time within a day ??

✓ How to refine the degenerate primers in homological cloning??

✓ …

tissue-specific

developmental
stage-specific

photoperiod

circadian clock

**ଔ  Before my starting to clone genes ...**

**Useful tools based on microarrays data**

◆ GENEVESTIGATOR

Estimate the tissue and developmental stage specificity

https://www.genevestigator.ethz.ch/

◆ Diurnal : only for Arabidopsis genes

Estimate diurnal and circadian gene expression profile

http://diurnal.cgrb.oregonstate.edu/

Example:  clone the *4CL1* gene in Arabidopsis

## ❧ Before my starting to clone genes …

ealy stage of flowering

seedling

ealy stage of fruiting



**Example**: clone the *4CL1* gene in Arabidopsis

What I used:

the whole plant of seedlings without tissue-bias

**Before my starting to clone genes …**

' the expression peak accurs at 4 hours after light '

Example:  clone the *4CL1* gene in Arabidopsis

ෆ **Before my starting to clone genes …**

Example:  clone the *4CL1* gene in Arabidopsis

All those information considered, the sample for mRNA isolation might be the young seedlings collected at 4 hours after the light turned on.



total RNA $\xrightarrow{\text{RT}}$ cDNA $\xrightarrow{\text{PCR}}$ 1686bp

" **Half day on the Web,**

   **saves you half month in the lab!** "

**Before my starting to clone genes ...**

How to refine the degenerate primers for homological cloning ??

✓ *4CL* gene cloning from swtichgrass

Primer-F

```
At4CL1 : DDNESVPIPEGCLRFTELTQSTTEA----SEVIDSVEISPDDVVALPYSSGTTGLPKGVMLTHKGIVTSVAQCVDGENENLYFHS-DDVI : 250
At4CL2 : DSD---AIPENCLRFSELTQSEEPR----VDSIP-EKISPEDVVALPFSSGTTGLPKGVMLTHKGIVTSVAQCVDGENENLYFNR-DDVI : 243
At4CL3 : DEP----TPENCLPFSTLITDDETN-----PFQETVDIGGDDAAALPFSSGTTGLPKGVVLTHKSLITSVAQCVDGDNENLYLKS-NDVI : 253
Ip4CL1 : DEDD--GTPDGCQPFWAIVSAADEN-----SVPESP-ISPDDAVALPYSSGTTGLPKGVVLTHGGIVSSVAQCVDGENENLHMRAGEDVV : 257
Ip4CL2 : DSA-----PDGCLHFSELTCADENE----APCVD---ISPDDVVALPYSSGTTGLPKGVMLTHKGLITSVAQCVDGDNENLYFHS-EDVI : 226
Ip4CL3 : DGR-----RDGCVDFAELIAGEELP-----EADEAGVLPDDVVALPYSSGTTGLPKGVMLTHRSIVTSVAQIVDGSNENVCFNK-DDAL : 237
Os4CL1 : DER-----RDGCLHFWDDIMSEDEASPIAGDEDDEKVFDPDDVVALPYSSGTTGLPKGVMLTHRSLSTSVAQCVDGENENIGLHA-GDVI : 249
Zm4CL  : DGR-----FDGCVEFAELIAAEEL-------EADADIHPDDVVALPYSSGTTGLPKGVMLTHRSLITSVAQCVDGENENLYFRK-DDVV : 243
```

```
At4CL1 : LCVLEMFHIYAINSIMLCGLRVGAAILIMPKFEINLLLELIQRCKVTVAFMVPPIVIAIAKSSETEKYDLSSIRVVKSGAAPLGKELEDA : 340
At4CL2 : LCVLEMFHIYAINSIMLCSLRVGATILIMPKFEITLLLEQIQRCKVTVAMVVPPIVIAIAKSPETEKYDLSSVRMVKSGAAPLGKELEDA : 333
At4CL3 : LCVLPLFHIYSINSVLINSLRSGATVLIMHKFEIGALLDLIQRHRVTIAAIVPPIVIAIAKNPTVNSYDLSSVRFVLSGAAPLGKELQDS : 343
Ip4CL1 : LCVLPLFHIFSINSVLLCALRAGAAVMIMPRFEMGAMLEGIERWRVTVAAVVPPIVIAIAKNPGVEKHDLSSIRIVLSGAAPLGKELEDA : 347
Ip4CL2 : LCVLEMFHIYAINSIMLCGLRVGAPILIMPKFEIGSLLGLIEKYKVSIAFVVPFVMMSIAKSPDLDKHDLSSLRMIKSGCAPLGKELEDT : 316
Ip4CL3 : LCLLPLFHIYSLHTVLIAGLRVGAAIVIMRKFDVGAIVDIVRAHRITIAPFVPPIVVEIAKSDRVGADDIASIRMVLSGAAEMGKDLQDA : 327
Os4CL1 : LCALEMFHIYSINTIMMCGLRVGAAIVVMRRFDIAAMMDIVERHRVTIAPIVPPIVVAVAKSEAAAARDLSSVRMVLSGAAEMGKDIEDA : 339
Zm4CL  : LCLLPLFHIYSINSVLIAGLRAGSTIVIMRKFDLGAIVDIVRRYVITIAPFVPPIVVEIAKSPRVTAGDIASIRMVMSGAAEMGKELQDA : 333
```

Primer-R

```
At4CL1 : VNAKFENAKLGQGYGMTEAGEVIAMSLGFAKEPFEVKSGACGTVVRNAEMKIVDPDTGDSLSRNQPGEICIRGHQIMKGYINNEAATAET : 430
At4CL2 : ISAKFENAKLGQGYGMTEAGEVIAMSLGFAKEPFEVKSGACGTVVRNAEMKILDPDTGDSLPRNKPGEICIRGNQIMKGYINDPIATAST : 423
At4CL3 : LRRRLPCAILGQGYGMTEAGEVLSMSLGFAKEPIPTKSGSCGTVVRNAELKVVHLETRLSLGYNQPGEICIRGQQIMKEYINDPEATSAT : 433
Ip4CL1 : LRGRLPCAIFGQGYGMTEAGEVLSMCEAFAREPTEAKSGSCGTVVRNAQLKVVDPDTGVSLGRNLPGEICIRGPQIMKGYINDFVATAAT : 437
Ip4CL2 : VRAKFPCARLGQGYGMTEAGEVIAMCLAFAKEPFDIKPGACGTVVRNAEMKIVDPETGASLPRNQPGEICIRGDQIMKGYINDPEATSRT : 406
Ip4CL3 : FMAKIFNAVLGQGYGMTEAGEVIAMCLAFAKEPFKVKSGSCGTVVRNAELKVVDPDTGASLGRNQPGEICVRGKQIMIGYINDPESTKNT : 417
Os4CL1 : FMAKLPGAVLGQGYGMTEAGEVLSMCIAFAKEPFKVKSGACGTVVRNAELKIIDPDTGKSLGRNLPGEICIRGQQIMKGYINNPEATKNT : 429
Zm4CL  : FMAKIFNAVLGQGYGMTEAGEVIAMCLAFAKEPYEVKSGSCGTVVRNAELKIVDPDTGAALGRNQPGEICIRGEQIMKGYINDPESTKNT : 423
```

## ∝ **Before my starting to clone genes ...**

How to refine the degenerate primers for homological cloning ??

✓ *4CL* gene cloning from swtichgrass

Degenerate primers:

SSGTTGLPKGV

**WSNWSNGGNACNACNGGNYTNCCNAARGGNGTN**

PGEICIRG

**CCNGGNGARATHTGYATHMGNGGN**

**So many uncertain sites, how to deal with ??**

## ∝ **Before my starting to clone genes ...**

✓ *4CL* gene cloning from swtichgrass

*What I did ...* **Condontree** for looking into condon bias



codontree : codon usage table, distance matrix and bases composition
(Pesole, Attimonelli and Liuni)

Reset | Run codontree | ● | [_____] your e-mail
(● = required, ● = conditionally required)

● Sequences File : please enter either :

1. the name of a file: [_____] 浏览...

2. *or* the **actual data** here:

(sequence format)

Control options

Output options

http://bioweb.pasteur.fr/seqanal/interfaces/codontree.html

## ✑ **Before my starting to clone genes …**

```
===========================================================
 Am. Acid    Codon        Number      Freq. %     Cod-Use
===========================================================
   Leu        CTA            5          0.37         0.05
   Leu        CTC           42          3.12         0.39
   Leu        CTG           43          3.20         0.40
   Leu        CTT            8          0.59         0.07
   Leu        TTA            5          0.37         0.05
   Leu        TTG            5          0.37         0.05

   Ser        AGC           19          1.41         0.18
   Ser        AGT            4          0.30         0.04
   Ser        TCA           16          1.19         0.16
   Ser        TCC           28          2.08         0.27
   Ser        TCG           19          1.41         0.18
   Ser        TCT           17          1.26         0.17

   Arg        AGA           14          1.04         0.17
   Arg        AGG           17          1.26         0.21
   Arg        CGA            4          0.30         0.05
   Arg        CGC           19          1.41         0.23
   Arg        CGG           22          1.64         0.27
   Arg        CGT            6          0.45         0.07

   Gly        GGA           11          0.82         0.11
   Gly        GGC           50          3.72         0.52
   Gly        GGG           25          1.86         0.26
   Gly        GGT           10          0.74         0.10

   Val        GTA            6          0.45         0.05
   Val        GTC           38          2.83         0.32
   Val        GTG           58          4.31         0.50
   Val        GTT           15          1.12         0.13
```
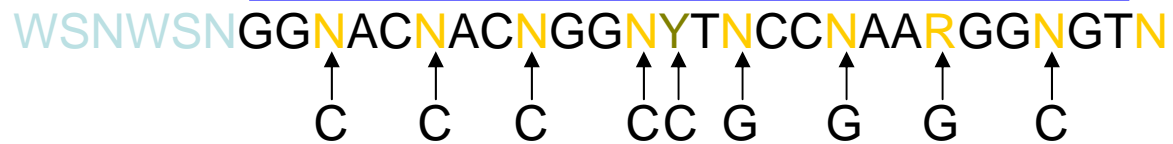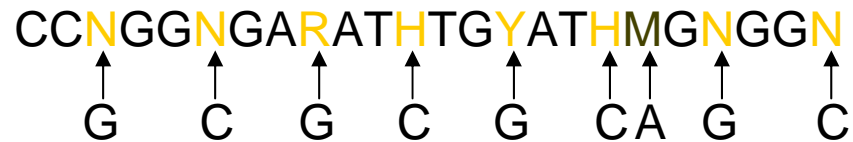
**…**

## ✀ **Before my starting to clone genes …**

✓ *4CL* gene cloning from swtichgrass

primer-F :  5'-CACCACCGGCCTGCCGAAGGGCGT-3'

WSNWSN GGNACNACNGGNYTNCCNAARGGNGTN
         C   C   C    CC G  G  G   C

primer-R :  5'-CGGGCGAGATCTGGATCAGGGG-3'

CCNGGNGARATHTGYATHMGNGGN
   G    C   G   C   G   CA G   C

# ❧ Bioinformatics in Plant Biology



*Annu. Rev. Plant Biol. 2006. 57:335–60*

# Acknowledgments

*Deeply appreciation to Prof. JC Luo for his precious edifications and instructions !*

*Cordial thanks to all my classmates who made the whole course study agreeable !*

Thanks for
your attentions !