

棉纤维伸长促进基因 *KCS* 的 生物信息研究

The bioinformatic research of
cotton fibre elongation-promoting
gene *KCS*





报告人：黄盖

组员：高晓雪、刘丹、屈颖

2014/1/10

Content



- 1 Background and information collection
 - 2 Genome-wide analysis
 - 3 Phylogeny construction and chromosome location
 - 4 3D structure prediction
 - 5 Summary
- 
- 

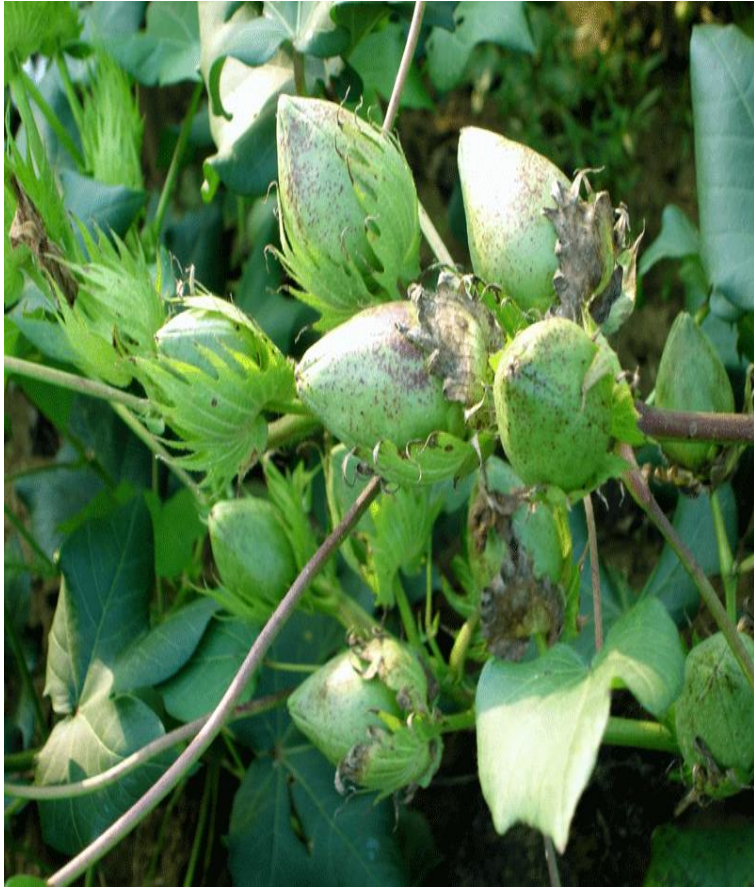




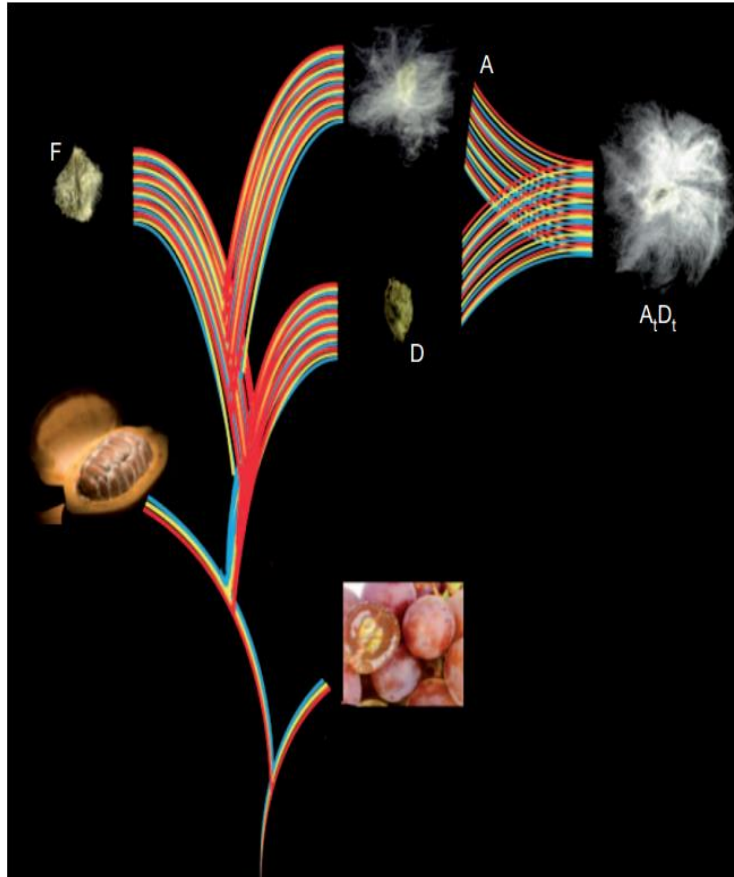
**Why we choose
cotton gene KCS?**



Cotton



About cotton evolution



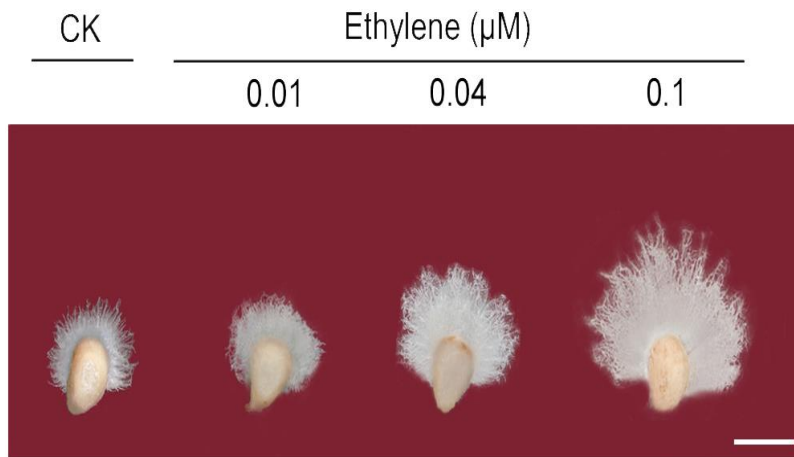
(Andrew H. et al.,2012)

- Cotton is an important economic crop.
- Gossypium genus contains tetraploid and diploid species. The tetraploid cotton species AD-genome (such as *G. hirsutum*) are thought to have formed by an allopolyploidization of A and D-genome species.
- *G. hirsutum* (AD-genome) produce most widely used natural fiber, but *G. raimondii* (D-genome) is fibreless.
- The draft genome of *G. raimondii* is sequenced by our lab.(Kunbo Wang. et al.,2012)

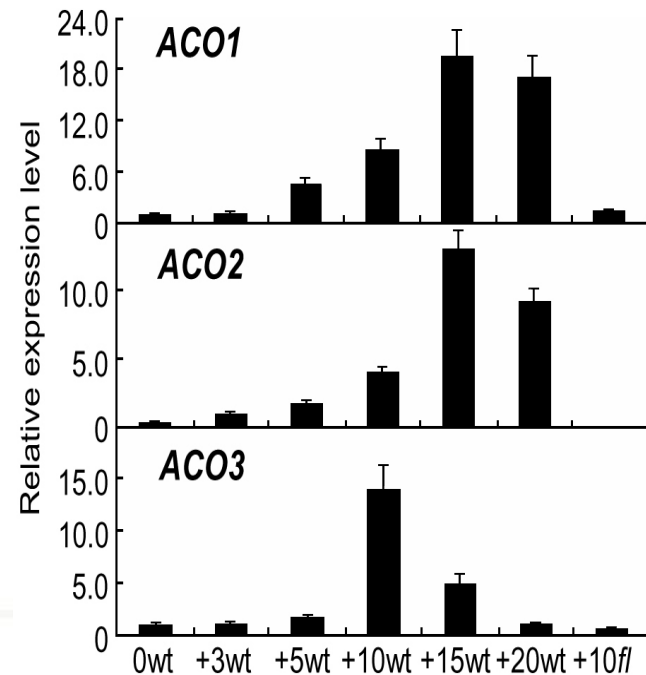
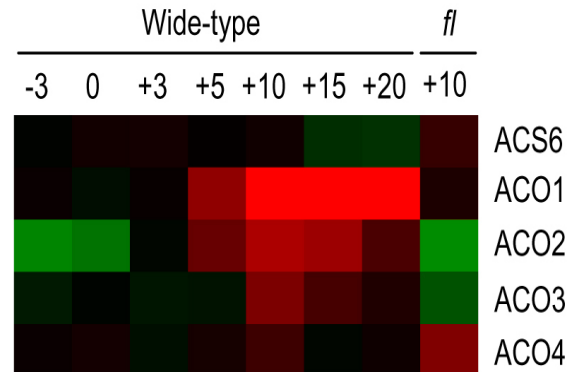


Previous Research of Our Lab

Ethylene

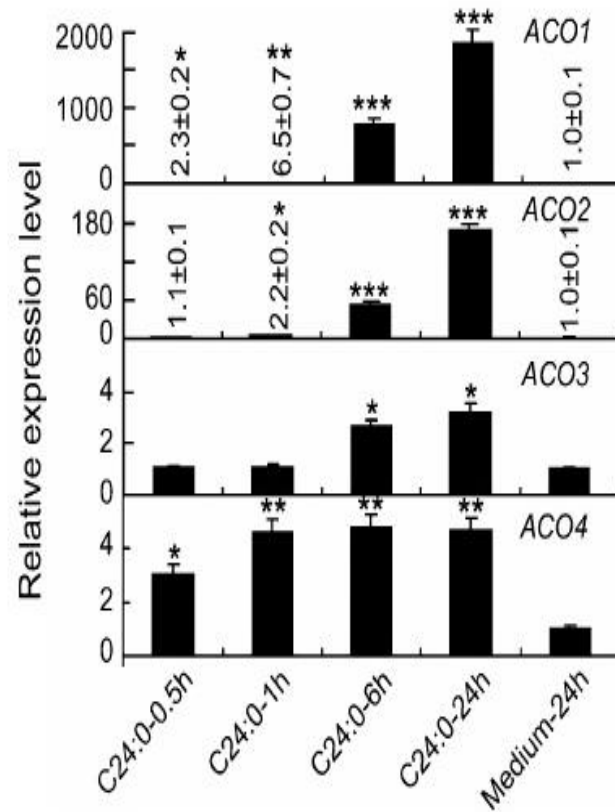
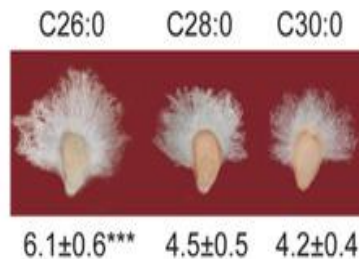
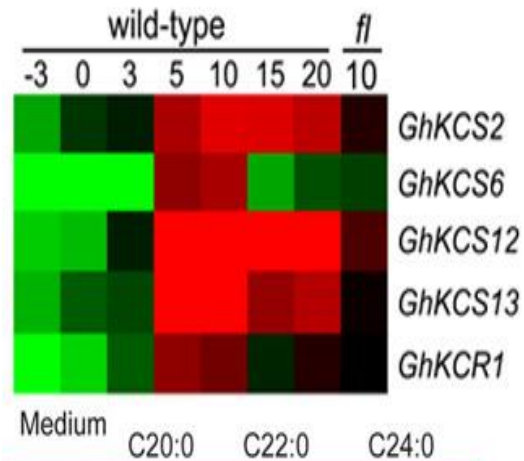


(Shi *et al.*, 2006)



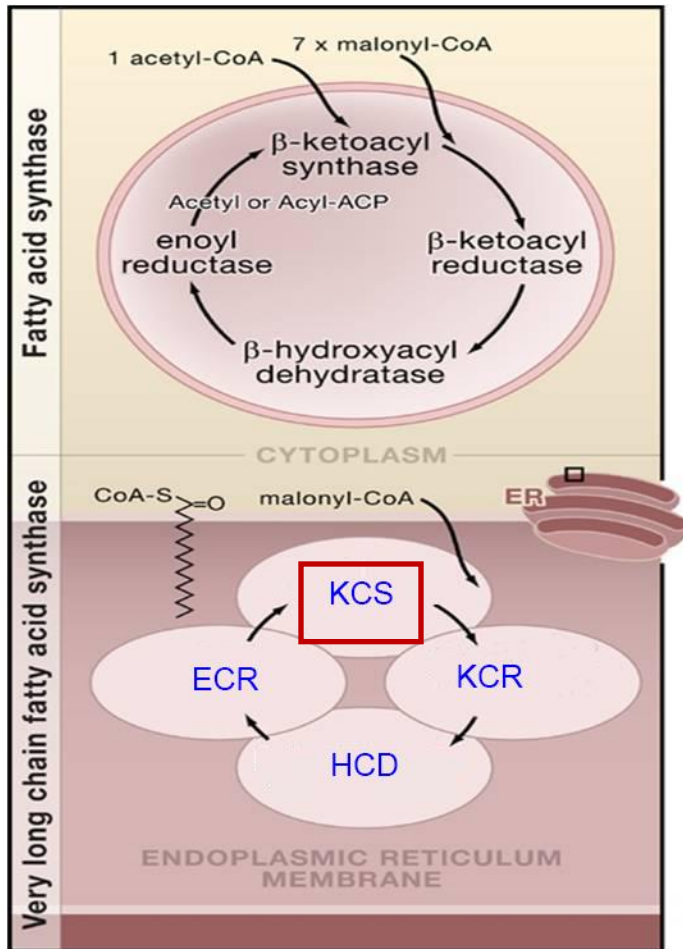
Previous Research of Our Lab

Very Long Chain Fatty Acids (VLCFAs)



(Qin et al., 2007)

About VLCFAs





(Riezman, 2007, revised)

- ◆ VLCFAs (>18 carbons are essential components of plant lipids, suberin and cuticular waxes.
- ◆ VLCFAs is synthesized by four successive enzyme reactions, including KCS, KCR, HCD and ECR
- Previous results in our lab reveal that VLCFAs may promote cotton fiber elongation by regulating ethylene biosynthesis.

UniProt database search

Gene names	Alias	Entry name	Length of amino acids
KCS1	EL1 At1g01120 T25K16.11	Q9MAM3_ARATH	528
KCS2	At4g34510 T4L20.90	O65677_ARATH	487
KCS3	At1g07720 F24B9.18	Q9LQP8_ARATH	478
KCS4	At1g19440 F18O14.21	Q9LN49_ARATH	516
KCS5	CER60 At1g25450 F2J7.9	Q9C6L5_ARATH	492
KCS6	CER6 EL6 At1g68530 T26J14.10	Q9XF43_ARATH	497
KCS7	At1g71160 F23N20.15	Q9C992_ARATH	460
KCS8	At2g15090 T15J14.13	Q4V3C9_ARATH	481
KCS9	At2g16280 F16F14.22	Q9SIX1_ARATH	512
KCS10	FDH EL4 At2g26250 T1D16.11	Q570B4_ARATH	550
KCS11	At2g26640 F18A8.1	O48780_ARATH	509
KCS12	At2g28630 T8O18.8	Q9SIB2_ARATH	476
KCS13	HIC At2g46720 T3A4.10	Q9ZUZ0_ARATH	466
KCS14	At3g10280 F14P13.12	Q9SS39_ARATH	459
KCS15	At3g52160 F4F15.270	Q9SU99_ARATH	451
KCS16	EL2 At4g34250 F10M10.20	Q9SYZ0_ARATH	493
KCS17	At1g04220 F20D22.1	Q5XEP9_ARATH	528
KCS18	FAE1 At4g34520 T4L20.100	Q38860_ARATH	506
KCS19	At5g43760 MQD19.11	Q9FG87_ARATH	529
KCS20	At5g49070 K20J1.4	Q9FH27_ARATH	464
KCS21	At5g04530 T32M21.130	Q9LZ72_ARATH	464

The database contains 21 reviewed sequences of Arabidopsis thaliana and they belong to FAE protein family. Besides, none of their protein has 3D structure .

Names	Alias	Entry name	Family	Length of AA
KCS2		A9XUG6_GOSHI		529
KCS1	Gr10021139	M9Z5H0_GOSRA		533
KCS2	Gr10031660	M9Z380_GOSRA		529
KCS3	Gr10022901	M9ZC32_GOSRA		466
KCS4	Gr10018297	M9Z385_GOSRA		504
KCS5	Gr10014500	M9ZC12_GOSRA		496
KCS6	Gr10032226	M9ZC26_GOSRA		467
KCS7	Gr10016173	M9ZA63_GOSRA		510
KCS8	Gr10026783	M9Z6N0_GOSRA		535
KCS9	Gr10018147	M9Z5I3_GOSRA		504
KCS10	Gr10015926	M9Z5H7_GOSRA		531
KCS11	Gr10017624	M9ZC02_GOSRA		510
KCS	Gr10018148	M9Z5I6_GOSRA		418
KCS13	Gr10031991	M9Z390_GOSRA		492
KCS15	Gr10032475	M9ZC21_GOSRA		462
KCS16	Gr10000033	M9Z6L3_GOSRA		390
KCS17	Gr10028536	M9ZC07_GOSRA		537
KCS18	Gr10034062	M9ZA66_GOSRA		515
KCS19	Gr10010162	M9Z6K4_GOSRA		512
KCS20	Gr10019136	M9Z3A9_GOSRA		437
KCS21	Gr10017854	M9Z397_GOSRA		457
KCS	Gr10019659	M9ZA81_GOSRA		439
KCS	Gr10004055	M9Z6K7_GOSRA		533
KCS	Gr10036653	M9ZC17_GOSRA		515
KCS	Gr10040229	M9Z5K7_GOSRA		501
KCS	Gr10018150	M9ZA87_GOSRA		493
KCS	Gr10018149	M9ZA72_GOSRA		504
KCS	Gr10040228	M9Z3B3_GOSRA		501
KCS	Gr10019657	M9Z6M6_GOSRA		424
KCS	Gr10019616	M9Z5J9_GOSRA		449
KCS	Gr10033317	M9Z3A3_GOSRA		274
KCS	Gr10018195	M9Z6L8_GOSRA		285
KCS14	Gr10009484	M9Z5J2_GOSRA		282
KCS12	Gr10018194	M9ZA77_GOSRA		286

The database contains 33 KCS sequences of *G. raimondii* and 1 KCS sequence of *G. hirsutum*. Besides, none of them are reviewed and their sequences were submitted by our lab.



enome-wide analysis



Cotton database construction

```
F:\blast\bin 的目录
2014/01/06 16:20 <DIR> .
2014/01/06 16:20 <DIR> ..
2013/03/13 06:34 5,488,640 blastdbcheck.exe
2013/03/13 06:34 6,592,512 blastdbcmd.exe
2013/03/13 06:34 4,174,336 blastdb_aliastool.exe
2013/03/13 06:34 9,608,704 blastn.exe
2013/03/13 06:34 9,604,608 blastp.exe
2014/01/06 19:12 88,207 blastpout
2013/03/13 06:34 9,597,440 blastx.exe
2013/03/13 06:34 9,463,808 blast_formatter.exe
2014/01/06 00:02 1,279 CDS.fas
2014/01/06 00:07 26,794 cdsout
2013/03/13 06:34 5,219,840 convert2blastmask.exe
2014/01/05 12:38 775,246,172 Cotton_D.all.anchored.fa
2014/01/05 19:03 337,026 Cotton_D.all.anchored.fa.nhr
2014/01/05 19:03 53,448 Cotton_D.all.anchored.fa.nin
2014/01/05 19:03 194,138,605 Cotton_D.all.anchored.fa.nsq
2013/03/13 06:34 9,780,224 deltablast.exe
2013/03/13 06:34 5,517,824 dustmasker.exe
2014/01/04 13:38 49,567,883 Gr-cds.fa.txt
2014/01/04 16:49 4,975,290 Gr-cds.fa.txt.nhr
2014/01/04 16:49 491,788 Gr-cds.fa.txt.nin
2014/01/04 16:49 11,348,776 Gr-cds.fa.txt.nsq
2014/01/04 17:07 18,193,161 Gr-PEP.fa.fasta
2014/01/04 17:08 5,016,266 Gr-PEP.fa.fasta.phr
2014/01/04 17:08 327,888 Gr-PEP.fa.fasta.pin
2014/01/04 17:08 15,083,278 Gr-PEP.fa.fasta.psq
2013/03/13 06:34 51,345 legacy_blast.pl
2013/03/13 06:34 6,111,232 makeblastdb.exe
2013/03/13 06:34 5,744,128 makenhindex.exe
2013/03/13 06:34 5,085,184 makeprofiledb.exe
2014/01/06 19:12 554 PEP.fas
2013/03/13 06:34 9,755,136 psiblast.exe
2013/03/13 06:34 9,621,504 rpsblast.exe
2013/03/13 06:34 9,613,312 rpstblastn.exe
2013/03/13 06:34 5,727,744 segmasker.exe
2013/03/13 06:34 9,716,224 tblastn.exe
```

1. Download local blast software from NCBI.
2. Download G. raimondii genome sequence from CGP.
3. To execute blast in DOS system

To acquire conserved domain



[HOME](#) | [SEARCH](#) | [BROWSE](#) | [FTP](#) | [HELP](#) | [ABOUT](#)



Pfam 25.0 (March 2011, 12273 families)

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. [More...](#)

QUICK LINKS

- [SEQUENCE SEARCH](#)
- [VIEW A PFAM FAMILY](#)
- [VIEW A CLAN](#)
- [VIEW A SEQUENCE](#)
- [VIEW A STRUCTURE](#)
- [KEYWORD SEARCH](#)
- [JUMP TO](#)

ANALYZE YOUR PROTEIN SEQUENCE FOR PFAM MATCHES

Paste your protein sequence here to find matching Pfam families.

This search will use an E-value of 1.0. You can set your own search parameters and perform a range of other searches [here](#).

Recent Pfam [blog](#) posts

Hide this

[No, seriously, we've made a release](#) (posted 1 April 2011)

Well, it should have been out about 6 months ago, but finally the long-awaited



pfam (<http://pfam.sanger.ac.uk/>)

Pfam search results

Significant Pfam-A Matches

[Show](#) or [hide](#) all alignments.

保守 domain 名称

Family	Description	Entry type	Clan
FAE1_CUT1_RppA	FAE1/Type III polyketide synthase-like p ...	Family	CL0046
#HMM	ylarrrrkvYLvdyacykpedelkvstetfleivkrvkkldesleflrkilersGlgteetyvPrsllleipeektlaeareEaeevlfgavdellaktkvkpkdigilvnc		
#MATCH	y+trtr tyLvdYacy p++lk+++ fte +k++ +de ++eFlrkItersGlgTet +P+st+ p++++atartEae v+fgatd lta ++vkp+digilvnc		
#PP	689*****		
#SEQ	VMTRKRFTYLVYACYLPPQNLKADWRFMEYAKEAADFDEPTMEFLRKIMERSGLGDETGAPPSMNCFFPFPMSMAARQEAELVMFGALDTLFASTNVKPRDIGILVNC		

复制保守 domain 序列

FAE1_CUT1-RppA HMM:

```
ylarrrrkvYLvdyacykpedelkvstetfleivkrvkkldesleflrkilersGlgteetyvPrsllleipeektlaeareEaeevlfgavdellaktkvkpkdigilvncslfsptPslsamvvnryklredvksynLsgmGCsaglisidlakdllqvhkntlalvvstEnitlnwYvGnersmllsnclFRvGgaavllsnksadrrrakykLvhhvRthkgaddkayrcvlqeeDeegkvGvslskdlvkvagealkknlttlgplvLPIsEklrflaslvarkl
```

To query cotton genome database

Length=290

Sequences producing significant alignments:

		Score (Bits)	E Value
Cotton_D_gene_10017624	locus=scaffold163:642071:643603:+	427	4e-147
Cotton_D_gene_10018297	locus=scaffold71:210359:211873:+	425	2e-146
Cotton_D_gene_10016173	locus=scaffold172:1629641:1631173:-	422	3e-145
Cotton_D_gene_10036653	locus=scaffold109:1190172:1191719:+	418	1e-143
Cotton_D_gene_10034062	locus=scaffold4:3957287:3958834:-	417	2e-143
Cotton_D_gene_10010162	locus=scaffold130:399751:401289:-	417	4e-143
Cotton_D_gene_10015926	locus=scaffold166:1016377:1017972:+	414	6e-142
Cotton_D_gene_10018147	locus=scaffold152:1567657:1569171:+	413	8e-142
Cotton_D_gene_10028536	locus=scaffold43:1655946:1657559:+	409	1e-139
Cotton_D_gene_10018150	locus=scaffold152:1584096:1585466:+	405	1e-139
Cotton_D_gene_10031660	locus=scaffold484:3835910:3838652:-	408	2e-139
Cotton_D_gene_10021139	locus=scaffold431:112479:114389:-	408	2e-139
Cotton_D_gene_10018149	locus=scaffold152:1577823:1579337:+	404	4e-138
Cotton_D_gene_10014500	locus=scaffold472:1357197:1358687:+	402	2e-137
Cotton_D_gene_10031931	locus=scaffold48:311495:312973:-	399	2e-136
Cotton_D_gene_1004465	locus=scaffold107:13904:139704:-	394	5e-132
Cotton_D_gene_10000033	locus=scaffold1232:972:2144:+	380	3e-130
Cotton_D_gene_10028783	locus=scaffold467:724594:726550:+	376	6e-127
Cotton_D_gene_10018148	locus=scaffold152:1574866:1576236:+	320	1e-106
Cotton_D_gene_10017854	locus=scaffold167:1150904:1152277:-	297	2e-097
Cotton_D_gene_10040228	locus=scaffold1:7996116:7997621:+	291	2e-094
Cotton_D_gene_10040229	locus=scaffold1:8009851:8011356:+	290	7e-094
Cotton_D_gene_10019136	locus=scaffold36:1496467:1497780:-	284	3e-092
Cotton_D_gene_10032226	locus=scaffold48:2307409:2308812:+	273	1e-087
Cotton_D_gene_10022901	locus=scaffold425:472748:474148:+	270	1e-086
Cotton_D_gene_10019616	locus=scaffold207:1262504:1263853:-	266	2e-085
Cotton_D_gene_10032475	locus=scaffold475:3526793:3528181:-	256	3e-081
Cotton_D_gene_1003743	locus=scaffold119:61298:691087:-	248	4e-081
Cotton_D_gene_10019659	locus=scaffold207:1692130:1693449:+	254	1e-080
Cotton_D_gene_10019657	locus=scaffold207:1657270:1658544:+	251	1e-079
Cotton_D_gene_10019658	locus=scaffold207:1664233:1664877:+	187	8e-058
Cotton_D_gene_10038985	locus=scaffold30:686293:686952:-	72.4	4e-015
Cotton_D_gene_10000813	locus=scaffold471:148816:149609:-	38.9	0.002
Cotton_D_gene_10026953	locus=scaffold260:135384:136656:-	39.3	0.003
Cotton_D_gene_10023921	locus=scaffold497:1006873:1009133:+	38.5	0.004
Cotton_D_gene_10002256	locus=scaffold414:316833:319857:-	37.4	0.009
Cotton_D_gene_10023920	locus=scaffold497:934662:936994:+	37.0	0.013
Cotton_D_gene_10029217	locus=scaffold496:1906837:1908121:-	37.0	0.014
Cotton_D_gene_10011223	locus=scaffold119:765522:768233:+	36.2	0.024
Cotton_D_gene_10007646	locus=scaffold236:85034:86759:-	35.0	0.059

1. To query cotton genome

database based on HMM model
sequence

2. To retrieve the corresponding
gene sequences ,then confirm the

conserved domain in **Pfam** database
and **smart database** .

(<http://smart.emblheidelberg.de/>)

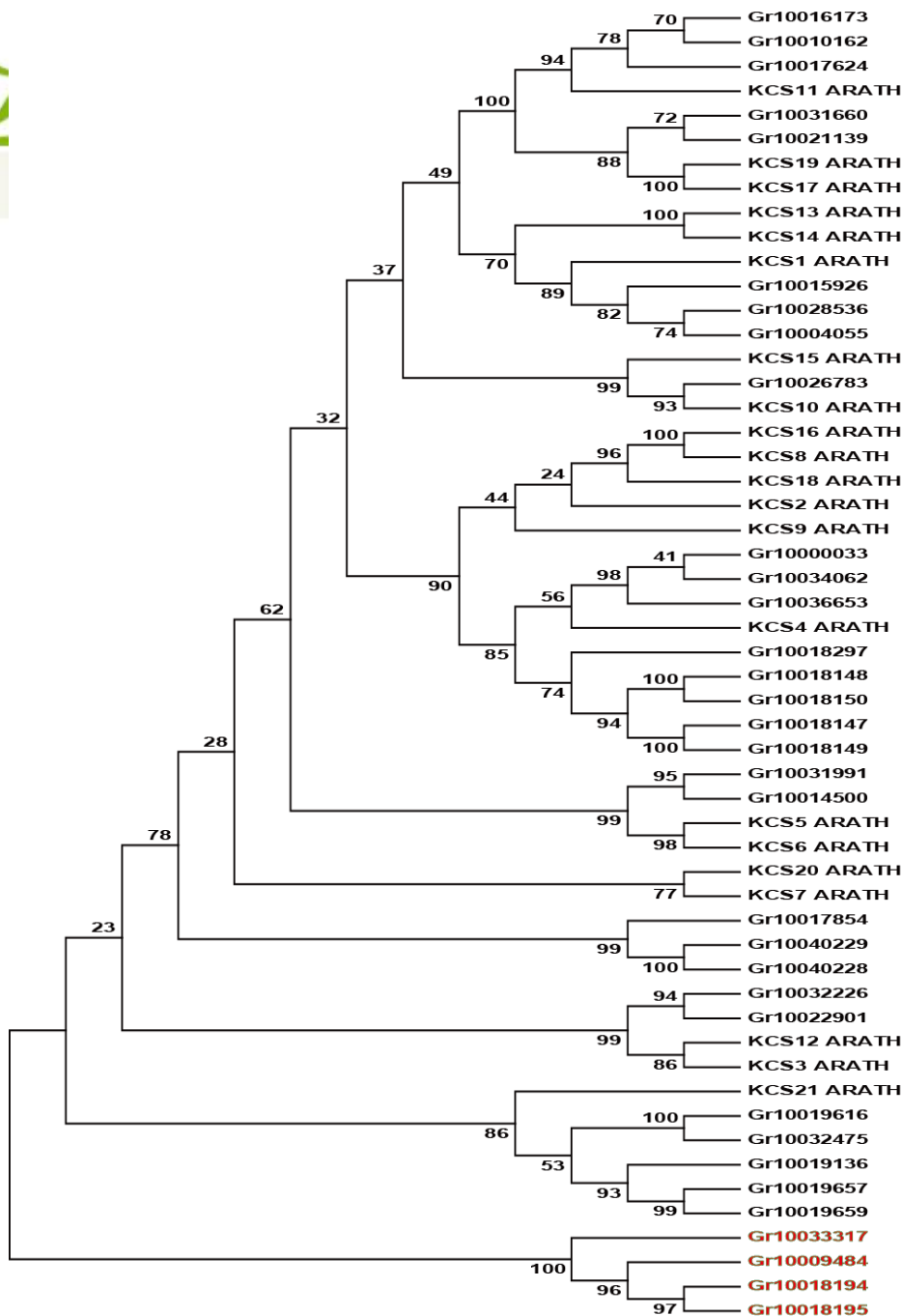
3. To further search the cotton

database in order to find all KCS
sequences.



Phylogeny construction and chromosome location





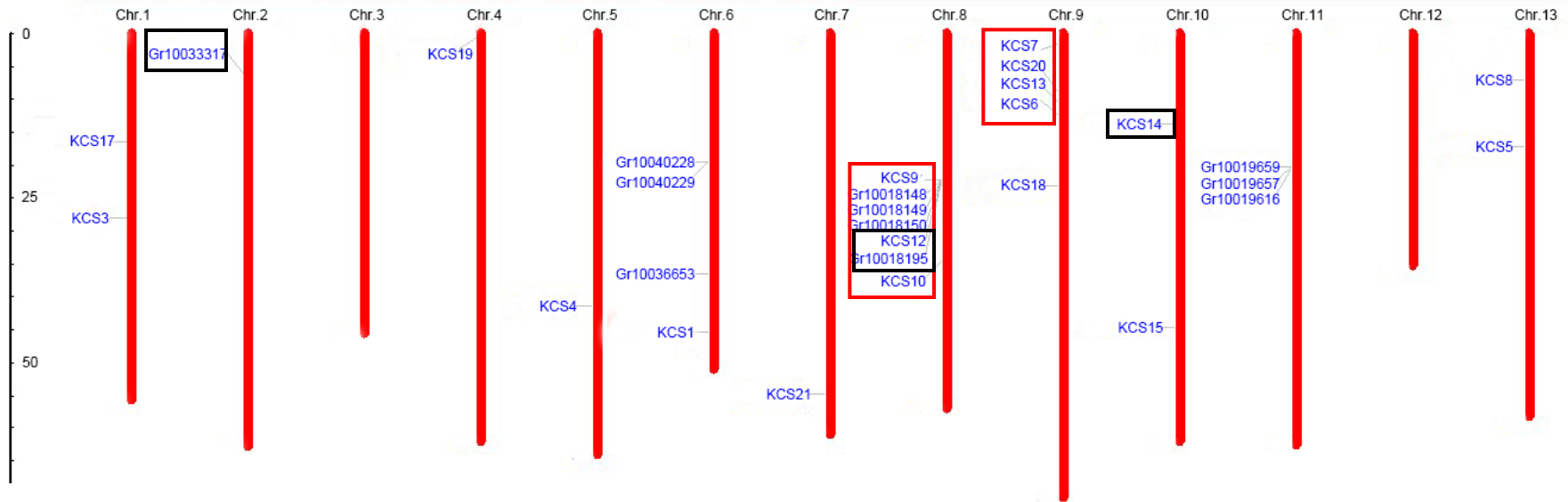
1. The phylogeny analysis suggests that all *G. ranondii* KCS genes are divided in two major branches.

2. In combination with Pfam prediction results, we know that the two major branches belong to **FAEFAE1_CUT1_RppA (PF08392 , PF08541)** and **Elo (PF01151, GNS1/SUR4 family)**

3. **FAE1_CUT1_RppA**: described as 3-ketoacyl-CoA synthases ; ACP_syn_III_C, ACP synthase III C terminal ; ELO1 Members of this family are involved in LCFAs elongation systems that produce the 26-carbon precursors for ceramide and sphingolipid synthesis

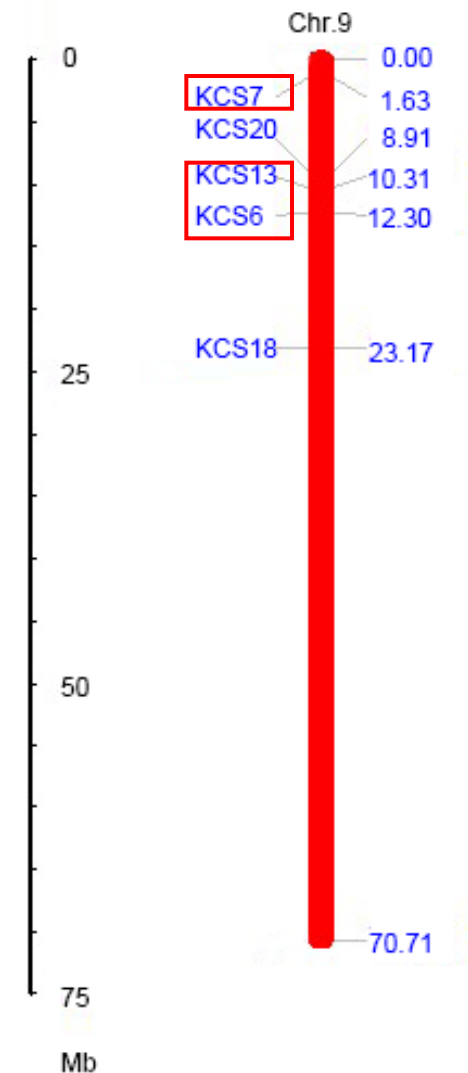
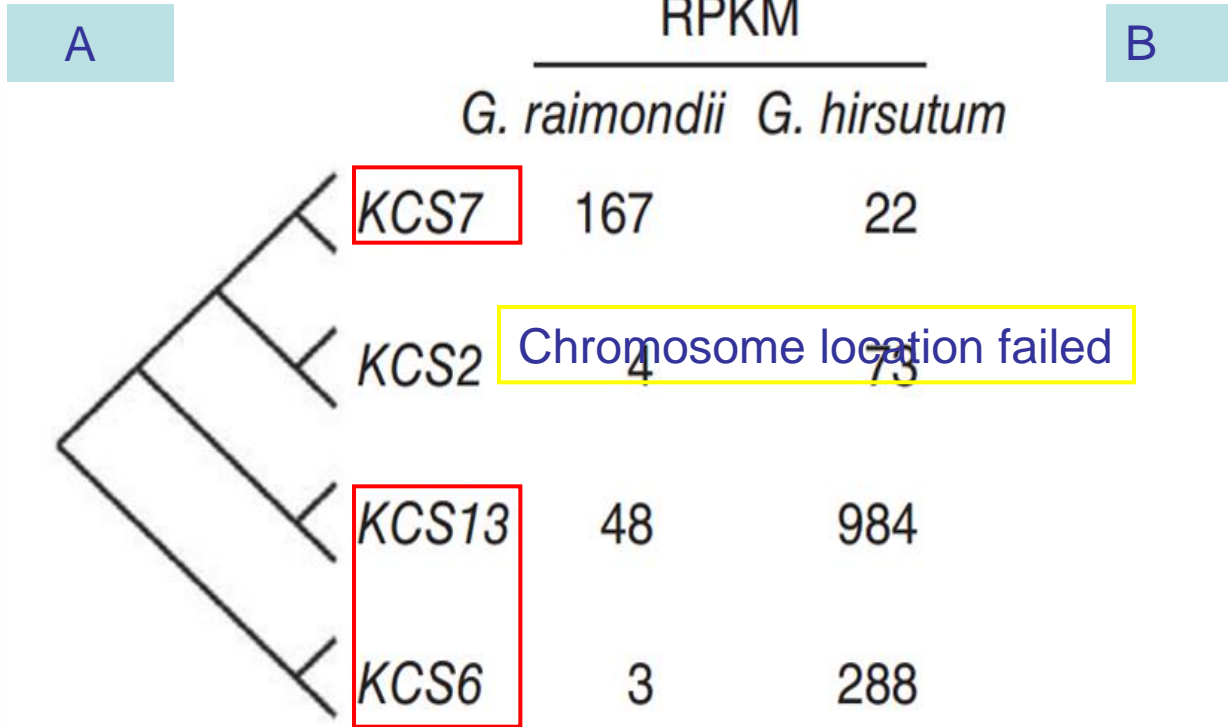


Candidate gene in chromosome location



29 *G. raimondii* KCS genes except GrKCS11 , GrKCS16, GrKCS2, Gr10004055 were located on 11 chromosomes ,but chromosome 3 and 12 contain no KCS genes.

KCS may be required in fiber initiation and elongation



Several 3-ketoacyl-CoA synthase (KCS) genes, including **KCS2**, **KCS13** and **KCS6**, were **only expressed in *G. hirsutum***, whereas intermediate levels of KCS7 transcripts were observed in both *G. hirsutum* and *G. raimondii*, indicating that high-level expression of Sus and KCS family genes may indeed be required for fiber cell initiation and elongation. (Kunbo Wang, et al. 2012. *Nature Genetics*.)

Multiple sequence alignment

```

380      390      400      410      420
Gr10040229 .....KEGEIYMPSSFRTAIQHFCLEPTSGRALIGETAKGLNDDGRDVEASLMTLRRF
Gr10040228 .....KEGEIYMPSSFRTAIQHFCLEPTSGRALIGETAKGLNDDGRDVEASLMTLRRF
Gr10017854 .....ESAEIYTPRFKTVVQHFCLEPSSGKPLIREVAKGLNENGRNIEPALMTLRRF
Gr10015926 .....AKVSPYIPDFKLAFDHFCHAGGRAVLDELQKNLQTDWHEPFRMTLRRF
Gr10028536 .....AKVSPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQTDWHEPFRMTLRRF
Gr10004055 .....ARVKPYIPDFKLAVEHFCIHAGGRAVLDELQKNLQTDWHEPFRMTLRRF
Gr10016173 .....MKVKPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQSEWHMPFRMTLRRF
Gr10017624 .....MKIKPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQSEWHMPFRMTLRRF
Gr10010162 .....MKIKPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQSEWHMPFRMTLRRF
Gr10031660 .....MKIKPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQTDWHEPFRMTLRRF
Gr10021139 .....MKIRPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQSDWHEPFRMTLRRF
Gr10031991 .....PKWKPYPDFKQAFEHFCIHAGGRAVLDELQKNLQSAEHVPSRMTLRRF
Gr10014500 .....PKWKPYPDFKLAFEHFCIHAGGRAVLDELQKNLQSAEHVPSRMTLRRF
Gr10036653 .....ASVKPYIPDFKLAFDHFCHAGGRAVLDELQKNLQLPIHVPASRMTLRRF
Gr10034062 .....AGIKPYIPDFKLAFDHFCHAGGRAVLDELQKNLQLPIHVPASRMTLRRF
Gr10000033 .....AGIKPYIPDFKLAFDHFCHAGGRAVLDELQKNLQLPIHVPASRMTLRRF
Gr10018297 .....AKIKPYIPDFKLAFDHFCHAGGRAVLDELQKNLQLPVHVPASRMTLRRF
Gr10018148 .....AKIKPYIPDFKLAFEHFCIHAGGRAVLDELQKNLQLPVHVPASRMTLRRF
Gr10018150 .....AKIKPYIPDFKLAFEHFCIHAGGRAVIGELKNLYLQVHVPASRMTLRRF
Gr10018147 .....AKIKPVVDFKLAFEHFCIHAGGRAVLDELQKNLQSEPLHVPASRMTLRRF
Gr10018149 .....AKIKPVVDFKLAFEHFCIHAGGRAVLDELQKNLQSEPLHVPASRMTLRRF
Gr10026783 KSKT...SLSPSSKPYIPDYKLAFEHFCVHAASKTVLDELQKNLQSENNMPASRMTLRRF
Gr10019616 .....KTSPPESGLNLSGGIDYFCIHPPGGRAVIDANGRLGNEYDLEPFRMALRRF
Gr10032475 .....NNNKGKTPSLNMXSGGFQHFCHHPGGRAVIDANGRLGNEYDLEPFRMALRRF
Gr10019657 .....KQGNLSFNLNLSGGVDHFCHHPGGRAVIDGLGKSLGSEYDLEPFRMALRRF
Gr10019659 .....KQGNLSFNLNLSGGVDHFCHHPGGRAVIDGLGKSLGSEYDLEPFRMALRRF
Gr10019136 .....KTAKP...SLNLSFGIQHFCHHPGGRAVIDGLGKSLGSEYDLEPFRMALRRF
Gr10032226 HGSSSHKGAQQGPIKAGVNFSGVDHFCHHTGGKAVIDGIGSLDLEYDLEPFRMALRRF
Gr10022901 HGSTKGTQTGGPIKAGVNFSGVDHFCHHTGGKAVIDGIGSLDLEYDLEPFRMALRRF

```

1. Align all KCS genes using ClustalX.
2. Resulting from poor alignment quality, we excluded all short sequences including Elo type protein and alignment again using ClustalX.

```

430      440      450      460      470      480
Gr10040229 GNSSSSMNYELAYNEAKREVRKKGDKVLMGLMCTGPKGSCVMECVRPPIAGDSNKNNEFR
Gr10040228 GNSSSSMNYELAYNEAKREVRKKGDKVLMGLMCTGPKGSCVMECVRPPIAGDSNKNNEFR
Gr10017854 GNSSSSMNYELAYNEAKREVRKKGDKIIVLGLCTGKCSLVLECLRPIVED.DKKSPPMS
Gr10015926 GNSSSSLNYELAYNEAKGRVSDGDRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10028536 GNSSSSLNYELAYNEAKGRVSSGDRVWQIAPGSOFPKNSAVNRRALRRSPMNELRGNPFX
Gr10004055 GNSSSSLNYELAYNEAKGRVSSGDRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10016173 GNSSSSLNYELAYNEAKGRERRGDRTWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10017624 GNSSSSLNYELAYNEAKGRERRGDRTWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10010162 GNSSSSLNYELAYNEAKGRIRKKGDRTWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10031660 GNSSSSLNYELAYNEAKGRERRGDRTWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10021139 GNSSSSLNYELAYNEAKGRIRKKGDRTWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10031991 GNSSSSLNYEMSYNEAKGRMKGDRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10014500 GNSSSSLNYEMSYNEAKGRMKGDRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10036653 GNSSSSLNYELAYNEAKGRIRNRRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10034062 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10000033 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10018297 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10018148 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10018150 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10018147 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10018149 GNSSSSLNYELAYNEAKGRMRRNRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10026783 GNSSSSLNYELAYNEAKREVRKKGDRVWQIAPGSOFPKNSAVNRRALRSTPMAESRGNPFX
Gr10019616 GNSAAGLYVLSVNEAKKRLKKGDRIMLSLCAQFPKNCNVCNVEVNMKDAMDDVNV...FK
Gr10032475 GNSAAGLYVLSVNEAKKRLKKGDRIMLSLCAQFPKNCNVCNVEVNMKDAMDDVNV...FK
Gr10019657 GNSAAGLYVLSVNEAKKRLKKGDKIFMVSLCAQFPKNCNVCNVEVNMKDGLDTRV...FE
Gr10019659 GNSAAGLYVLSVNEAKKRLKKGDKIFMVSLCAQFPKNCNVCNVEVNMKDGLDTRV...FE
Gr10019136 GNSAAGLYVLSVNEAKKRLKKGDKIFMVSLCAQFPKNCNVCNVEVNMKDGLDTRV...FE
Gr10032226 GNSAASLLNYVLAYNEAKKRLKKGDKVLMISLCAQFPKNSCLWVEVRRD.LGDGNV...FK
Gr10022901 GNSAASLLNYVLAYNEAKKRLKKGDKVLMISLCAQFPKNSCLWVEVRRD.LGDGNV...FK

```

3. The alignment results suggest that FAE type KCS protein share some strict conserved amino acids.
4. Based on Natural Evolution of Kimura, most amino acids mutants are neutral, few mutants are harmful and these sites are functional.

3D structure prediction



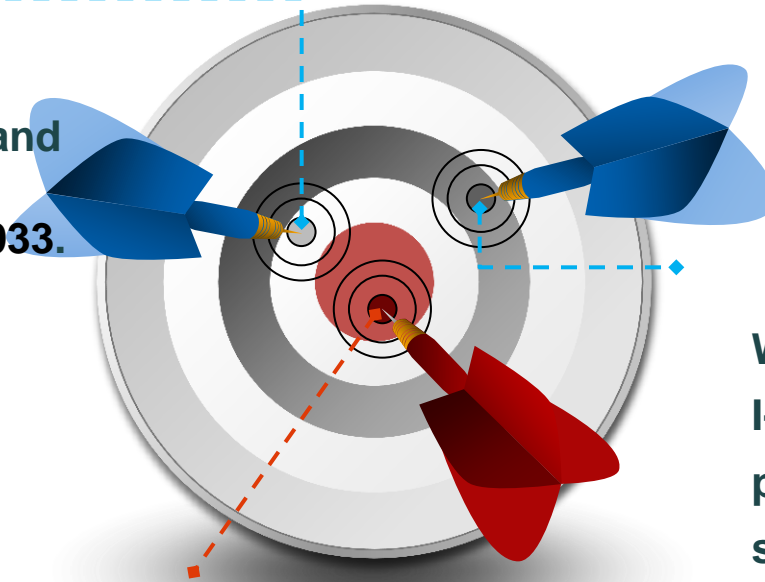
In order to find KCS proteins' active sites and study functional mechanism, we used popular methods to predict our target proteins' structures.



3D structure prediction

Obtain Squence

We extracted each one sequence from FAE type and Elo type , such as Gr10018148 and Gr10000033.



Prediction

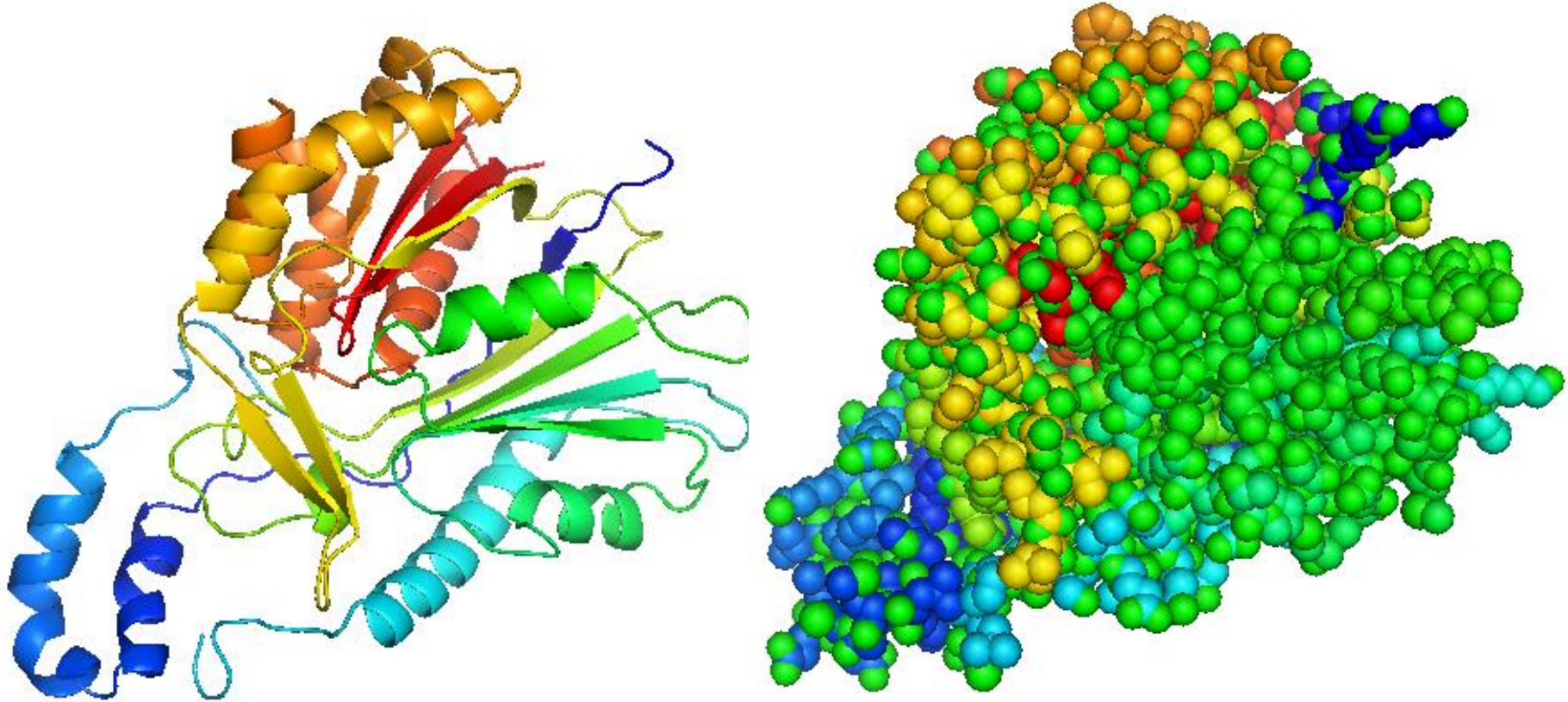
We used Phyre and I-TASSER methods to predict KCS protein 3D sturcuture.

Analysis

To analysis some details about protein structure



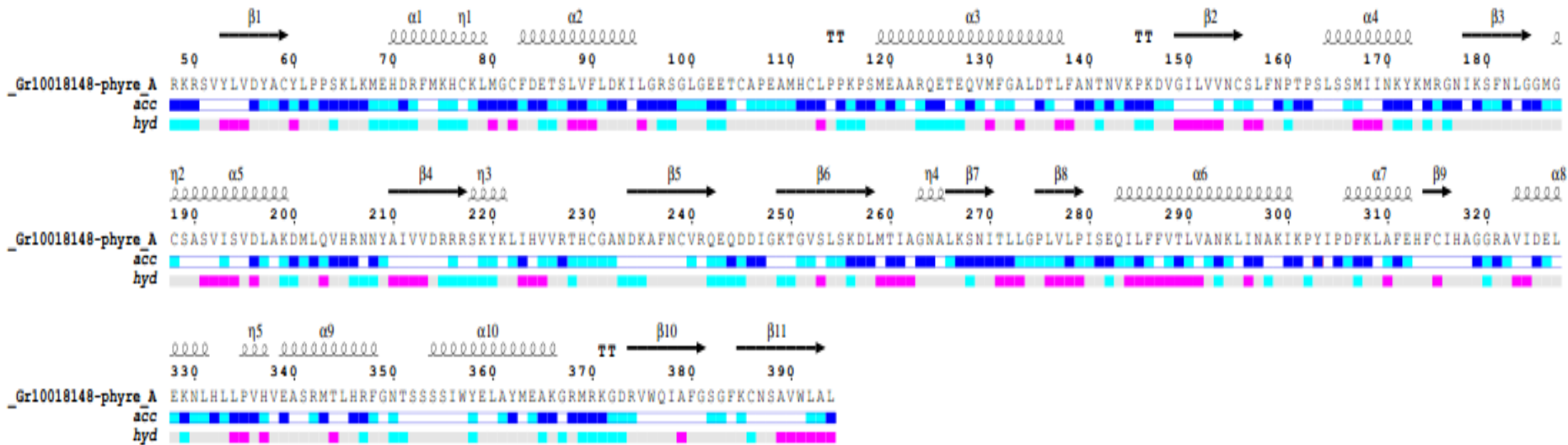
Analysis of Gr10018148 Protein structure



Gr10018148-phyre



The secondary structure of Gr10018148 protein



Acc: accessibility, white is buried; cyan is intermediate; blue is accessible; blue with red borders is highly exposed.

Hyd: hydrophobic, pink is hydrophobic; grey is intermediate; cyan is hydrophilic

Conserved domain in all KCS proteins

Gr10018148

160 α4 β3 η2 α5 210

Gr10018148 FNPT PSL S MI INK Y KMR GN I K S FN L G GMGCS A S V I S VDLAK DML QVHR N NYA I VV

Gr10040229 FCPS PSL S I I V NK Y SMR SD V K S FN L S GMGCS A G A I G IDLAQN LL KTNN N CYA I V I S T E I

Gr10040228 FCPS PSL S I I V NK Y SMR SD V K S FN L S GMGCS S G A I G IDLAQN LL KTNN N CYA I V I S T E I

Gr10017854 FCPS PSL S I I INK Y SMK SD I K S YN L S GMGCS A G T I G VDLAQN LL K T H E N K T A I V L S T E I

Gr10015926 FNPT PSL S A M I V N H Y K L R T N I N S YN L G GMGCS A G L I S IDLAKN LL Q S N P N T Y A L V V S T E N

Gr10028536 FNPT PSL S A M I V K H Y K L R T D I K S YN L G GMGCS A G L I S V E L A K N LL Q A N P N T Y A V V V S T E N

Gr10004055 FNPT PSL S A M V I N H Y K L R T D I K S YN L G GMGCS A G L I S V E L A K N LL R A N P N T Y A V V V S T E N

Gr10016173 FNPT PSL S A M V I N H Y K L R G N I Q S YN L G GMGCS A G L I S I D L A K H LL Q V H P N S Y A L V I S M E N

Gr10017624 FNPT PSL S A M V I N H Y K L R G N I Q S YN L G GMGCS A G L L S I D L A K N LL Q V H P N S Y A L V I S M E N

Gr10010162 FNPT PSL S A M V I N H Y K L R G N I Q S YN L G GMGCS A G L I S I D L A K N LL Q V H P N T Y A L V I S M E N

Gr10031660 FNPT PSL S A V V V N R Y K F R G N I L S YN L G GMGCS A G L I S I D L A K Q LL R V H P N S Y A L V V S M E N

Gr10021139 FNPT PSL S A M I V N R Y K L R G N I L S YN L G GMGCS A G L I S I D L A K Q M L Q V H P N S Y A L V V S M E N

Gr10031991 FSPT PSL S A M V I N K Y K L R S N I K S F N L S GMGCS A G L I S I D L A R D LL Q V H P N S N A V V V S T E I

Gr10014500 FSPT PSL S A M V I N K Y K L R S N I K S F N L S GMGCS A G L I S I D L A R D LL Q V H P N S N A V V V S T E I

Gr10036653 FNPT PSL S A M I V N K Y K L R G N I R S F N L G GMGCS A G V I A V D L A K D LL Q V H R N T Y A V V V S T E N

Gr10034062 FNPT PSL S A M I V N K Y K L R G N I R S F N L G GMGCS A G V I A V D L A K D M L Q V H R N T Y A V V V S T E N

Gr10000033 FNPT PSL S A M I V N K Y K L R G N I R S F N L G GMGCS A G V I A V D L A K D M L Q V H R N T Y A V V V S T E N

Gr10018297 FNPT PSL S A M I I N K Y K L R G N I R S F N L G GMGCS A G V I A I D L A K D M L Q V H R N S Y A V V V S T E N

Gr10018150 FNPT PSL S S M I I N K Y K M R G N I K S F N L G GMGCS A S V I S V D L A K D M L Q V H R N N Y A I V V S T E N

Gr10018147 FNPT PSL T A M I I N K Y K M R G N I K S F N L S GMGCS A G V I A I D L A K D M L Q V H R N N Y A V V F S T E N

Gr10018149 FNPT PSL T A M I I N K Y K M R G N N K S F N L S GMGCS A S V I A I D L A K D M L Q V Y R N N Y A V V F S T E N

Gr10026783 FNPT PSL S A M I I N H Y K M R G N I L S YN L G GMGCS A G I I A V D L A R D M L Q A N P N N Y A V V V S S E M

Gr10019616 I S S V P S I P A R V I N R Y K M R E D V K V F N L S GMGCS A S L I A V D L V N H L F Q T Y K N Q F A I V V S S E S

Gr10032475 I T S V P S L P A R V I N R Y K M R D D V K V F N L S GMGCS A S V I A V D L V H H L F K T Y M N S F A V I V S S E S

Gr10019657 F S P A P S L T S R I I N R Y K M R H N I K S F S L S GMGCS A S M L A I D L V Q Q L F K T Y K N Q F A I V V S T E S

Gr10019659 F S P A P S L T S R I I N R Y K M R D N I K S F S L S GMGCS A S M V A I D M V Q Q L F K T Y K N Q F A I V V S T E S

Gr10019136 F S P S P S L T A R I V N R Y K M R D N I K S F S L S GMGCS A S M V A I D L V Q N L F K S Y K N A F A V V V S S E T

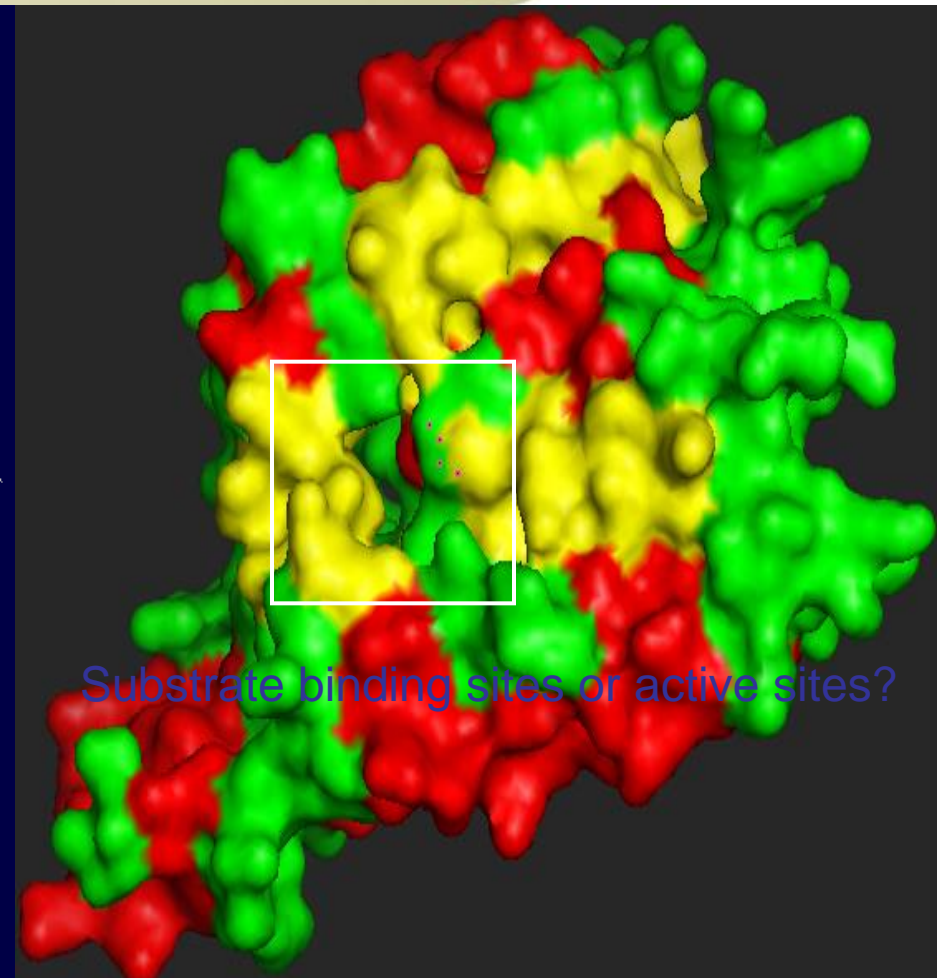
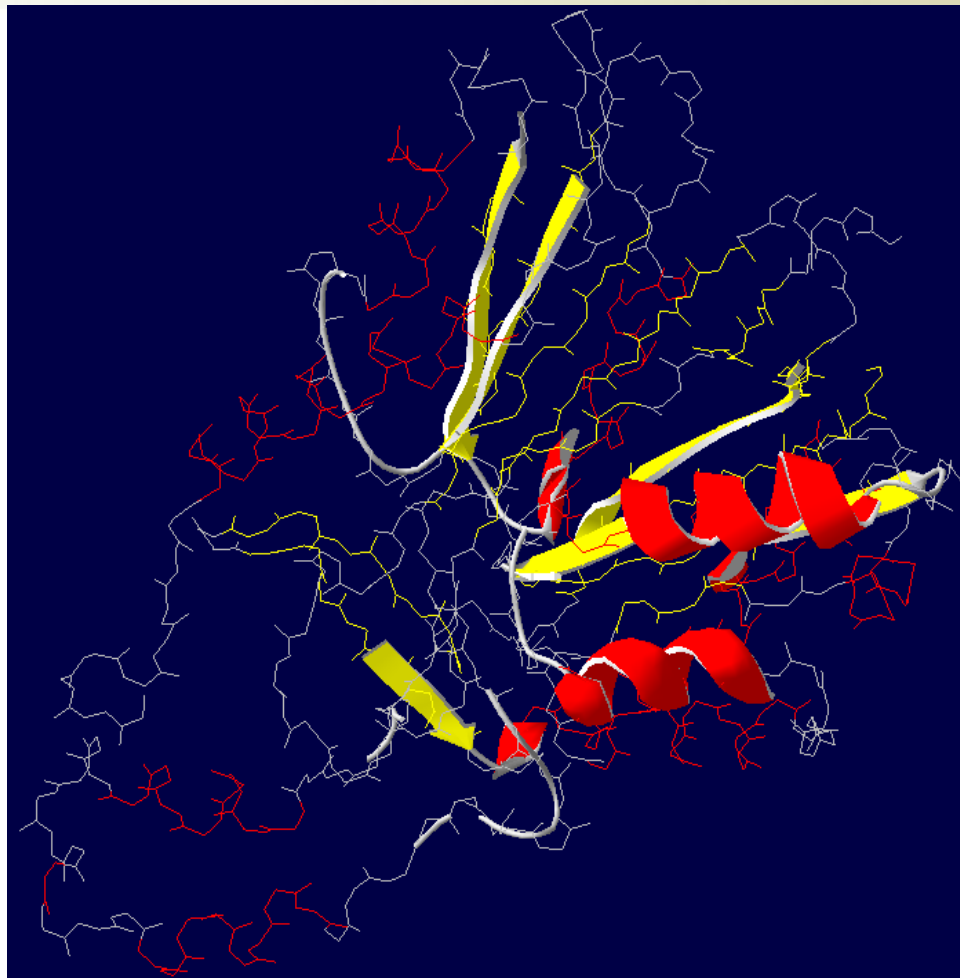
Gr10032226 L S T V P S L C S R I I N H Y K M R P D I K S F N L T GMGCS A S L I S L D I V R N V F K S Y K N K F A L L V T S E S

Gr10022901 I T A P P C L S S R I I N H Y K M R Q D I K C F N L T GMGCS A S L I S L D I V R N V F K S Y K N K Y A L L V T S E S

consensus>70 f . p . P s l s . . ! ! n . Y k m r . # i . s % n L . G M G C S a . . i . i # l a . d l l N . y A v v v s . e .



Why these domains are conserved?



Substrate binding sites or active sites?

Red: helix Yellow: sheet Loop: green

3D structure prediction suggest that Strict domains locate on two helix ,five sheets and some loops.



Summary

- ★ Local cotton database construction is helpful to query target sequence and conserved domain prediction contributes to research loss of function in protein family.
- ★ Phylogeny construction and chromosome location suggest that there are some KCS genes that are required for fiber initiation and elongation.
- ★ Protein structure prediction can provide useful information to design point mutants.
- ★ Bioinformatics analysis is useful, but we should be careful.



Acknowledgement

- Prof. Luo
- Prof. Zhu
- All members in ABC class, especially G01 group.
- Our lab members

