



Sequencing transcriptomes
in toto

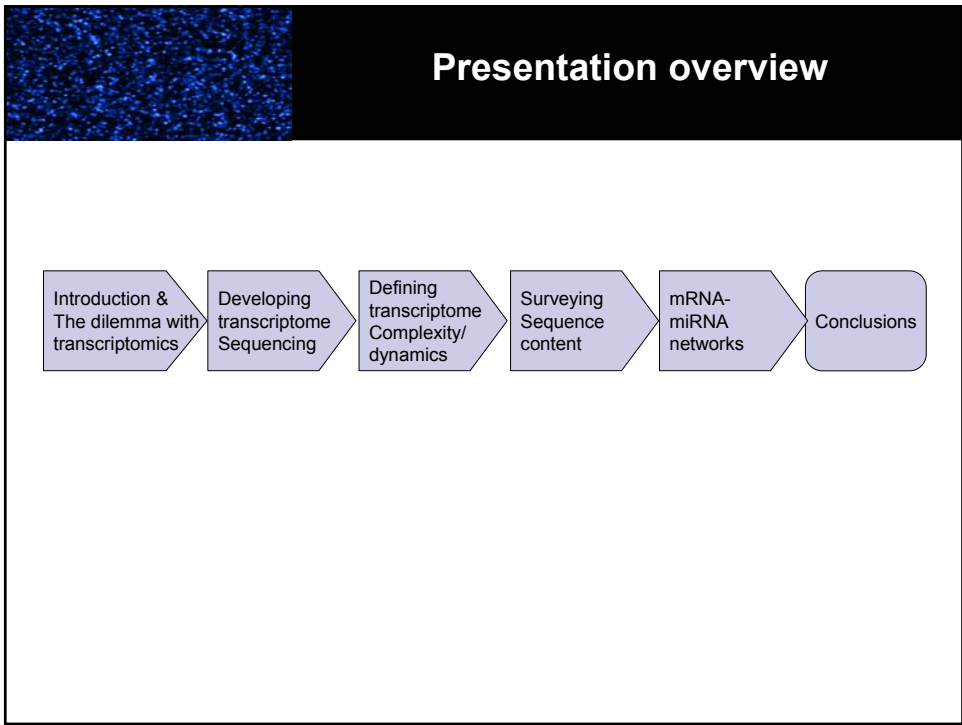
Sean Grimmond
November 12th 2008



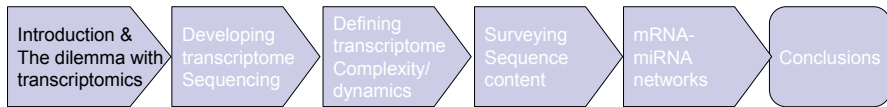
THE UNIVERSITY
OF QUEENSLAND
AUSTRALIA



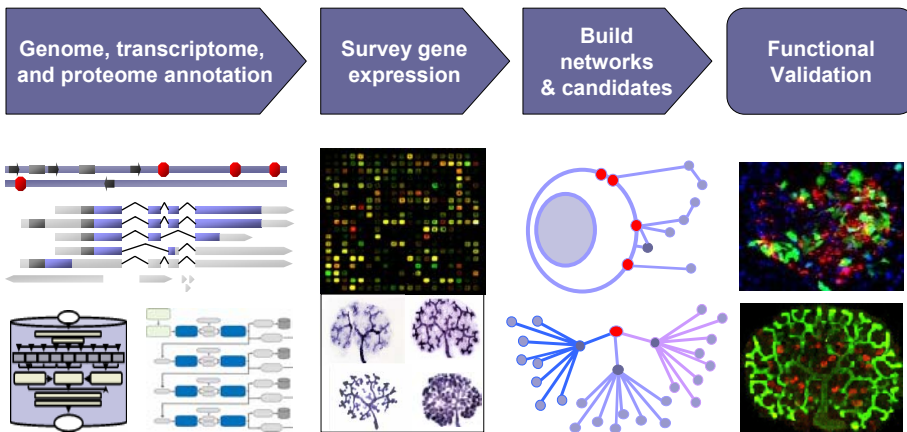
IMB Institute for Molecular Bioscience



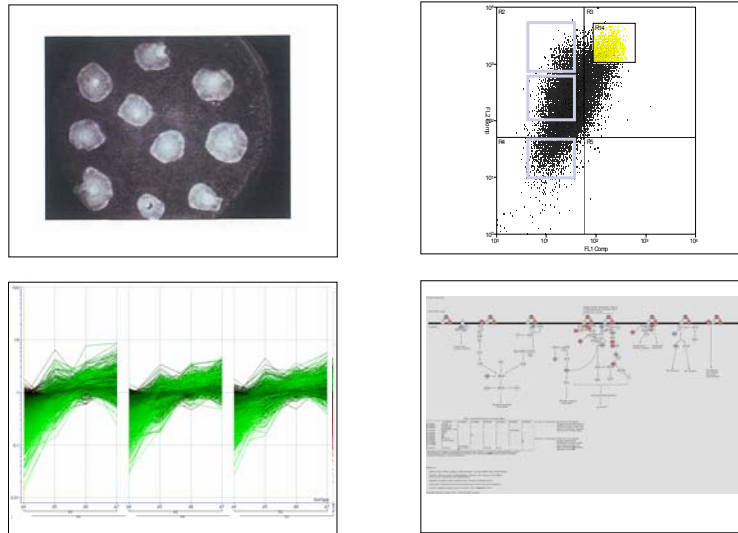
Presentation overview



Deriving biological insight from the transcriptome

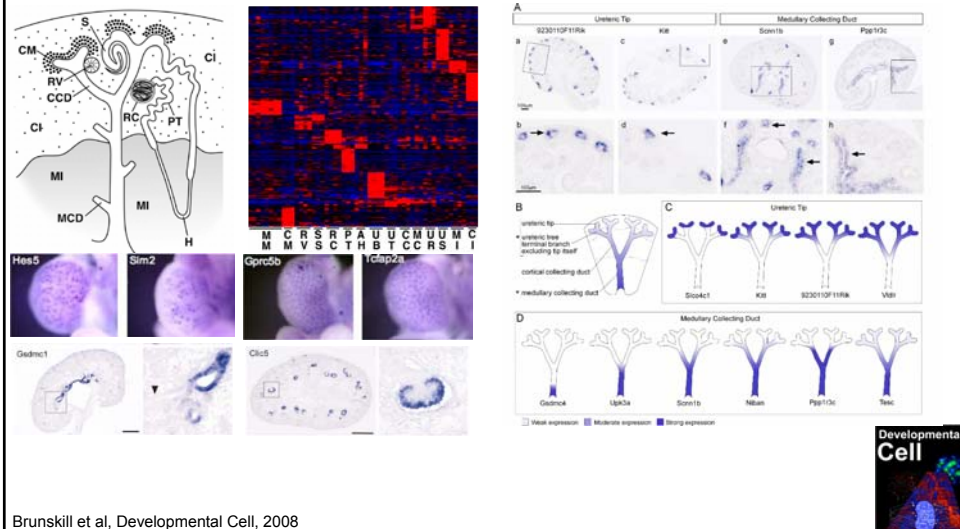


Transcriptome based screening of Human ES cell surface markers



Laslett et al, BMC Dev Biol. 2007

In situ based profiling of Urogenital development



Brunskill et al, Developmental Cell, 2008

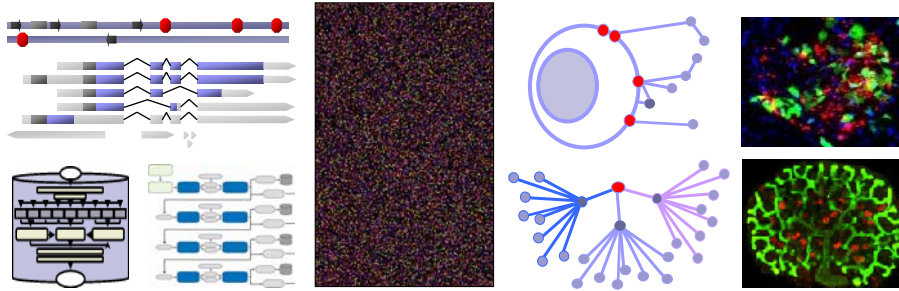
Deriving biological insight from the transcriptome

Genome, transcriptome, and proteome annotation

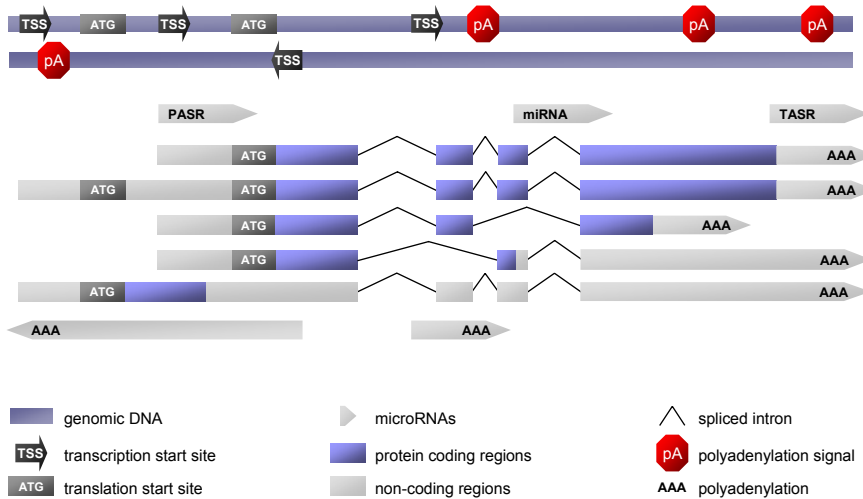
Survey mRNA
Polysome mRNA
miR levels

Build networks,
find candidates

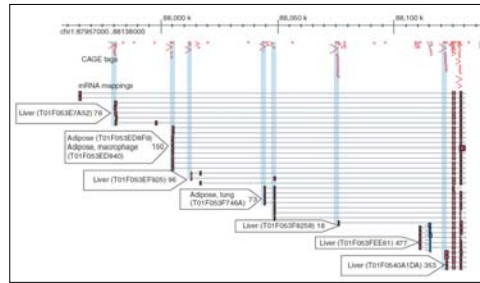
Functional
Validation



Representative Gene Vs Transcriptional complexity



Creating transcriptional frameworks against the genome:

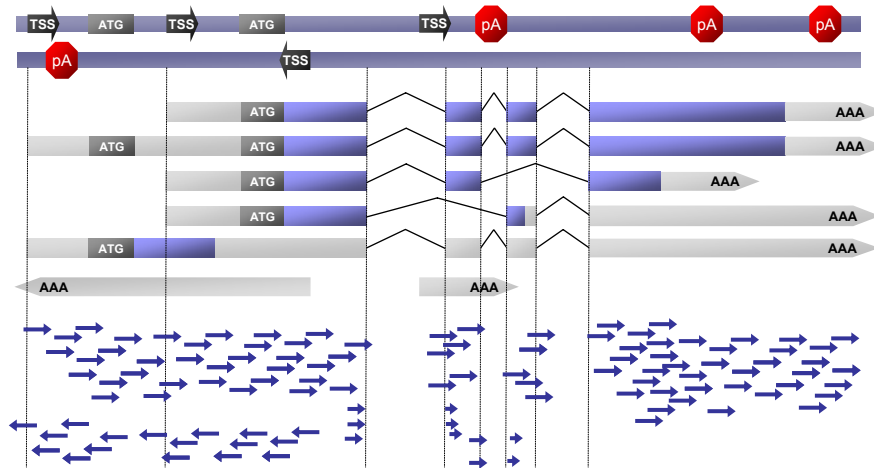


**Exhaustive surveying of expression space:
 ~30,000 genes, >100,000 transcripts, > 65,000 ORFs
 Multiple promoters, Novel layers of control...**

Carninci et al, Nat Genet, 2005.



Using shotgun sequencing to survey transcriptomes?



Presentation overview



Sequencing by ligation

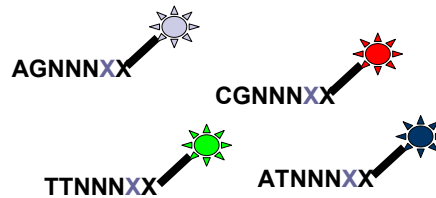
Accurate Multiplex Polony Sequencing of an Evolved Bacterial Genome

Jay Shendure,^{1*} Gregory J. Porreca,^{1*} Nikos B. Reppas,¹
Xiaoxia Lin,¹ John P. McCutcheon,^{2,3} Abraham M. Rosenbaum,¹
Michael D. Wang,¹ Kun Zhang,¹ Robi D. Mitra,² George M. Church¹

We describe a DNA sequencing technology in which a commonly available, inexpensive epifluorescence microscope is converted to rapid nonelectrophoretic DNA sequencing automation. We apply this technology to resequence an evolved strain of *Escherichia coli* at less than one error per million consensus bases. A cell-free, mate-paired library provided single DNA molecules that were amplified in parallel to 1-micrometer beads by emulsion polymerase chain reaction. Millions of beads were immobilized in a polyacrylamide gel and subjected to automated cycles of sequencing by ligation and four-color imaging. Cost per base was roughly one-ninth as much as that of conventional sequencing. Our protocols were implemented with off-the-shelf instrumentation and reagents.

Shendure et al (2005) Science 309:1728

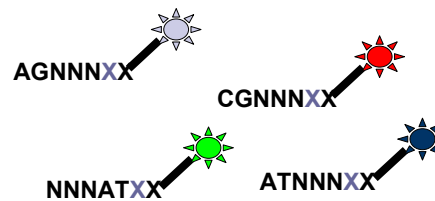
SOLiD sequencing



Sequencing by Supported Oligo Ligation and Detection (SOLiD).

1. Anneal sequencing primer to amplicon.
2. Anneal and ligate first dinucleotide detection probe. From pool of 16.
3. Wash and capture image of slide.
4. Cleave after dinucleotide detector and wash.
5. Repeat cycle 5-7 times (25-35 base read).

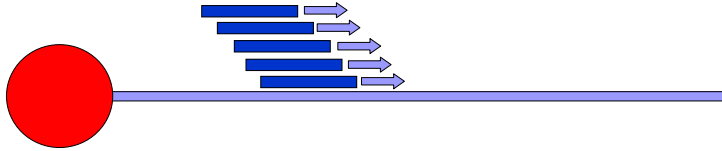
SOLiD sequencing



Sequencing by Supported Oligo Ligation and Detection (SOLiD).

1. Anneal sequencing primer to amplicon.
2. Anneal and ligate first dinucleotide detection probe. From pool of 16.
3. Wash and capture image of slide.
4. Cleave after dinucleotide detector and wash.
5. Repeat cycle 5-7 times (25-35 base read).

SOLiD sequencing



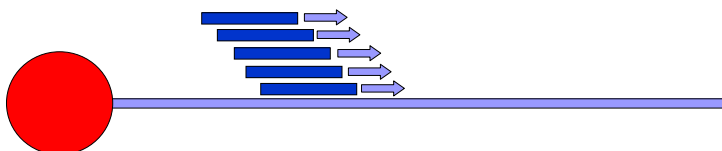
```

PrimerA:  T T 03 04 05  A T 08 09 10  C C 13 14 15  G G 18 19 20 ..
PrimerB:  T 02 03 04  C A 07 08 09  G C 12 13 14  T G 17 18 19  A ..
PrimerC:  01 02 03  A C 06 07 08  A G 11 12 13  A T 16 17 18  T A ..
PrimerD:  01 02  G A 05 06 07  A A 10 11 12  A A 15 16 17  C T 20 ..
PrimerE:  01  T G 04 05 06  T A 09 10 11 12 13 14 15 16  G C 19 20 ..
    
```

Sequence T T G A C A T A A G C C A A T G G C T A

- No homo-polymer problems, no phase problems
- Errors are random
- Each base is _read_twice_
- Independent measures of each base >> raw accuracy (99.5% for single runs)
- Consensus accuracy is 99.97% at 10x (based on human BAC sequencing)
- Output = 80-160,000,000 reads = 6.0Gb-18.0 Gb per run (3-9Gb per slide)
- Latest output is > 200,000,000 reads per slide at 50mers per run.

SOLiD sequencing



```

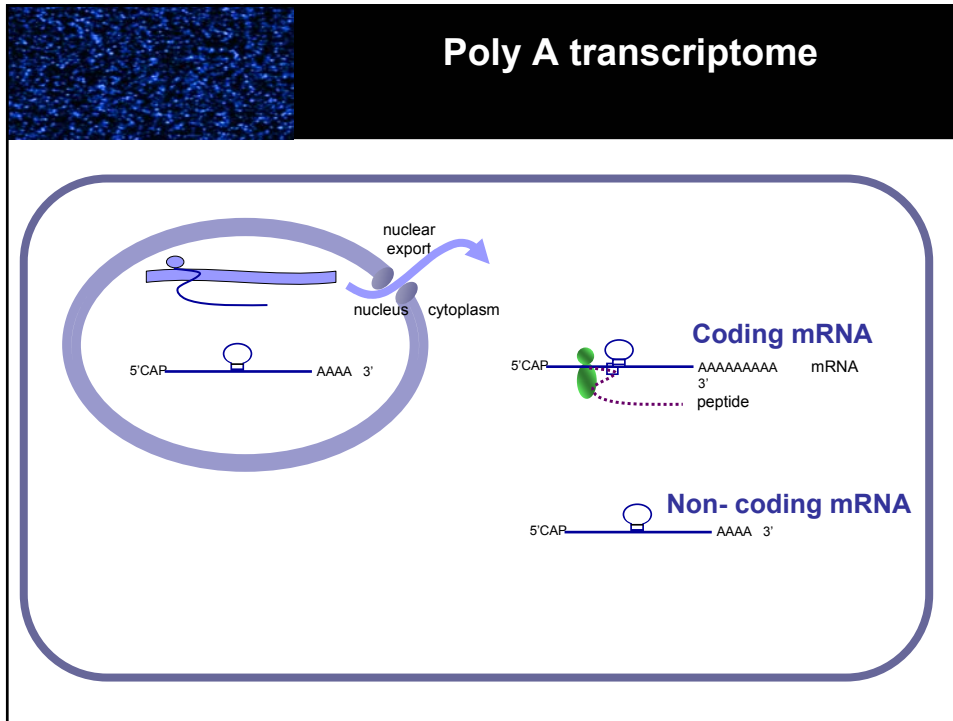
PrimerA:  T T 03 04 05  A T 08 09 10  C C 13 14 15  G G 18 19 20 ..
PrimerB:  T 02 03 04  C A 07 08 09  G C 12 13 14  T G 17 18 19  A ..
PrimerC:  01 02 03  A C 06 07 08  A G 11 12 13  A T 16 17 18  T A ..
PrimerD:  01 02  G A 05 06 07  A A 10 11 12  A A 15 16 17  C T 20 ..
PrimerE:  01  T G 04 05 06  T A 09 10 11 12 13 14 15 16  G C 19 20 ..
    
```

Sequence T T G A C A T A A G C C A A T G G C T A

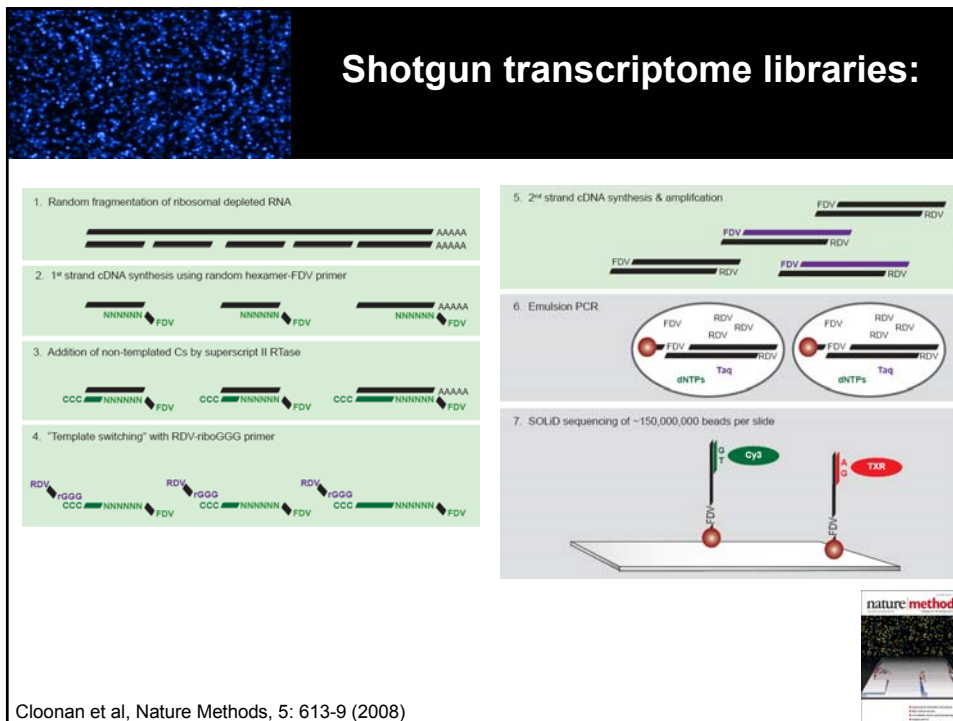
Scale:

- Output = 80-160,000,000 reads = 6.0Gb-18.0 Gb per run (3-9Gb per slide)
- Latest output is > 200,000,000 reads per slide at 50mers per run.

Poly A transcriptome

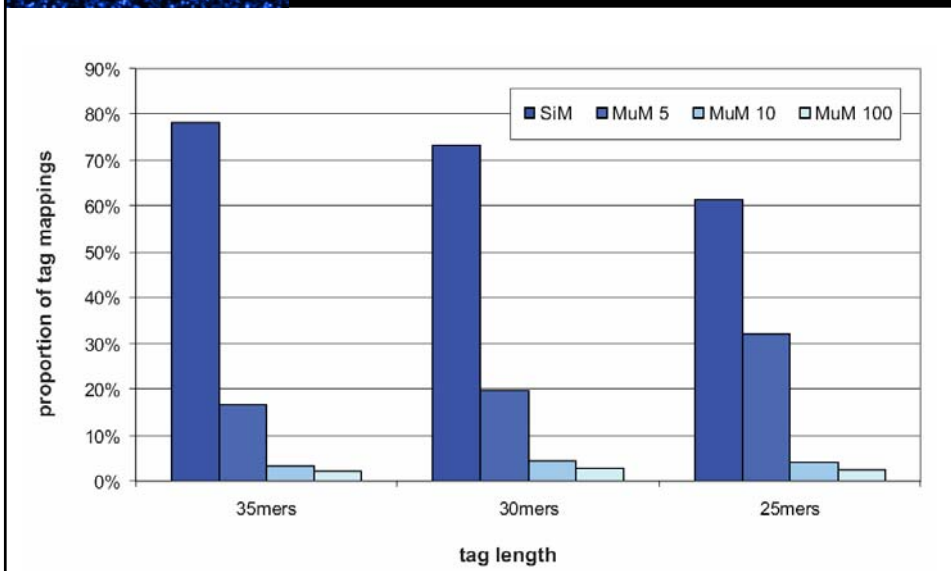


Shotgun transcriptome libraries:

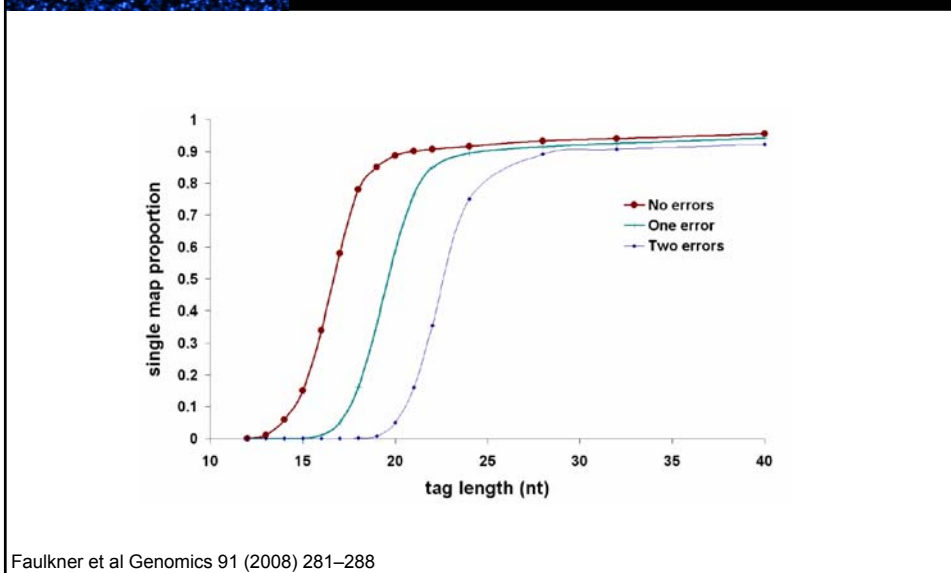


Cloonan et al, Nature Methods, 5: 613-9 (2008)

Accuracy of mapping at 25-35bp

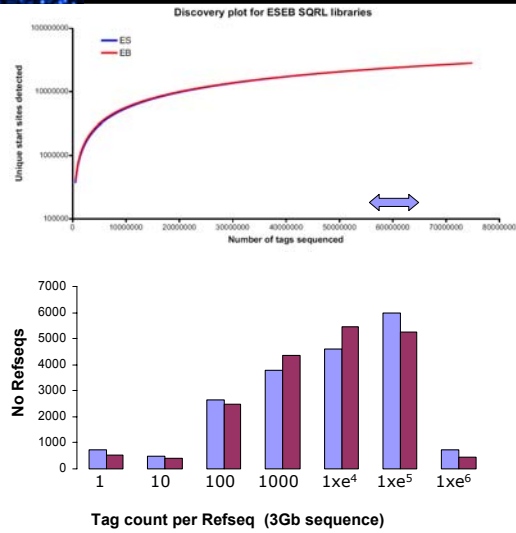


Unique mapping of short tags (impact of length & error):

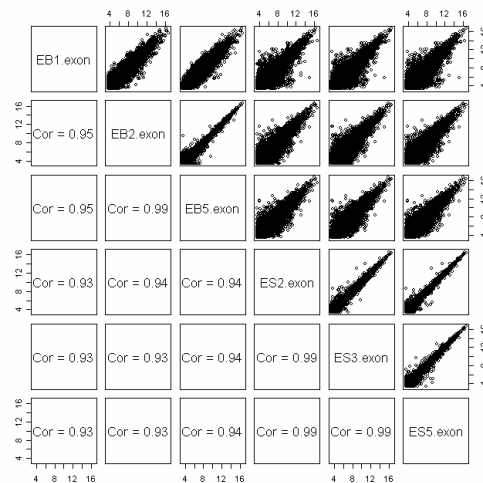


Faulkner et al Genomics 91 (2008) 281–288

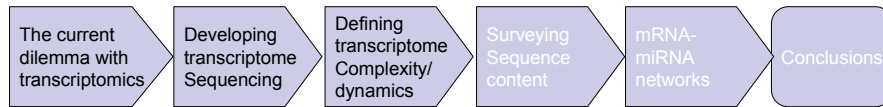
Depth of sequencing required for SQRL of mRNA



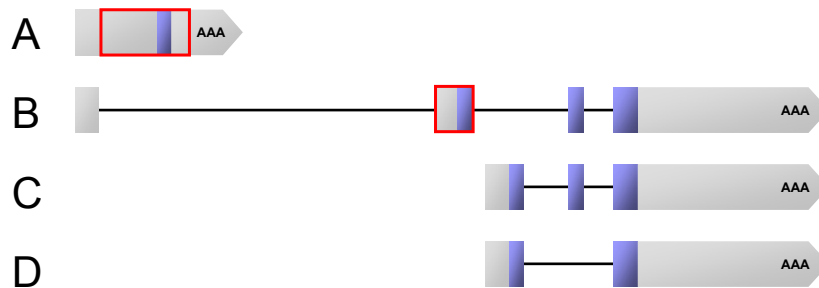
Biological and technical replication of SQRL



Presentation overview



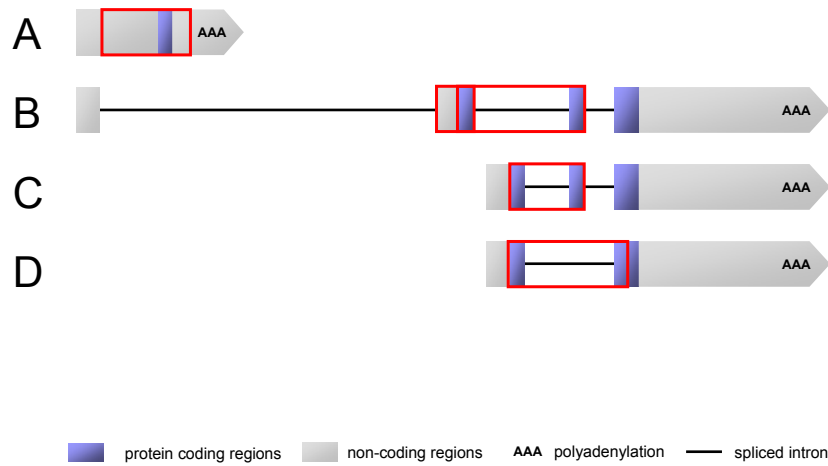
Defining transcript Specific expression



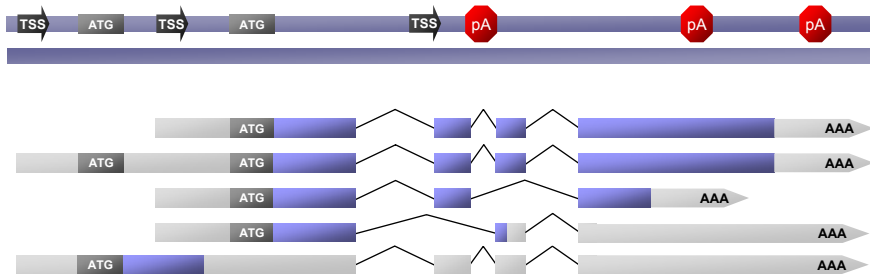
Transcripts defined by Aceview (September 2007 release)

■ protein coding regions ■ non-coding regions AAA polyadenylation — spliced intron

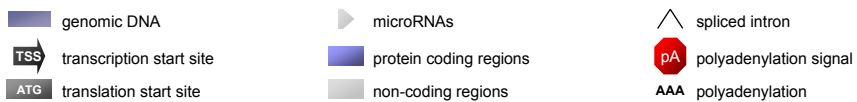
“Diagnostic” features



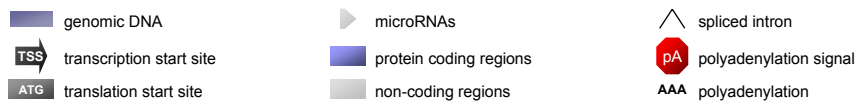
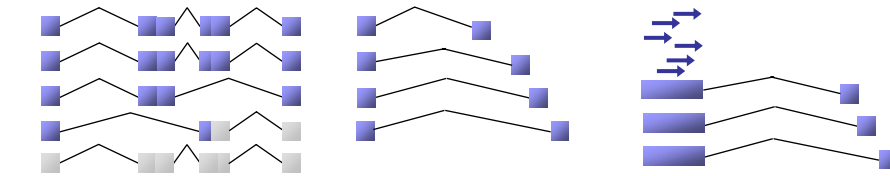
Defining transcript Specific expression



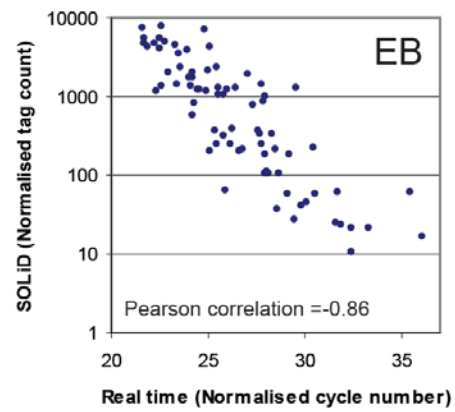
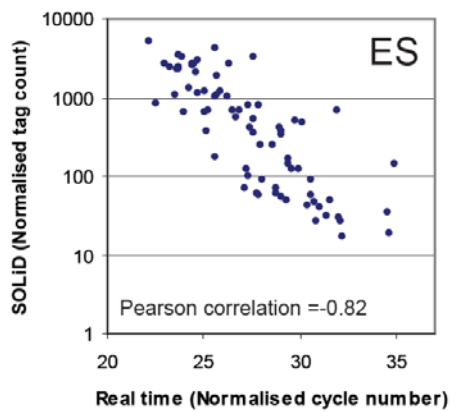
92.6% known transcripts have diagnostic features (covers 99.8% of loci)
 217127 diagnostic features covering 160156 individual transcripts from 65254 loci



Surveying exon usage



SQRL is quantitative correlation with qRT-PCR



Key questions in ES cell biology

Pluripotent Stem Cells

What are the signals required to maintain stem cell Pluripotency?

What are the molecular programs controlling lineage determination?

What role does RNA play in ES cell maintenance and differentiation?

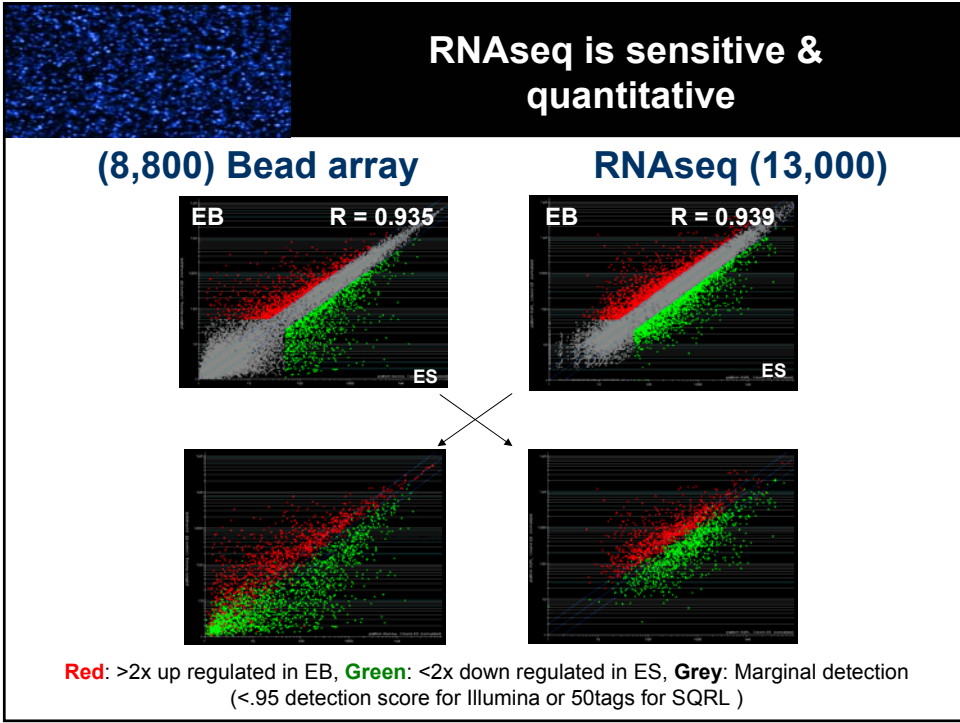
<http://www.stemcellresearchfoundation.org/>

Using shotgun sequencing to survey transcriptomes

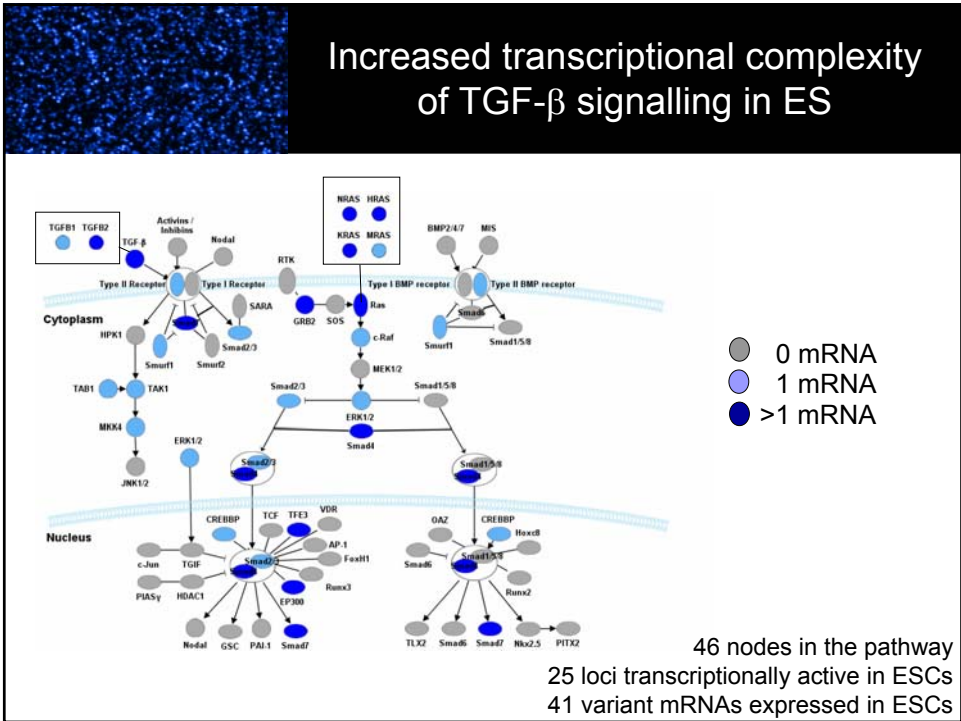
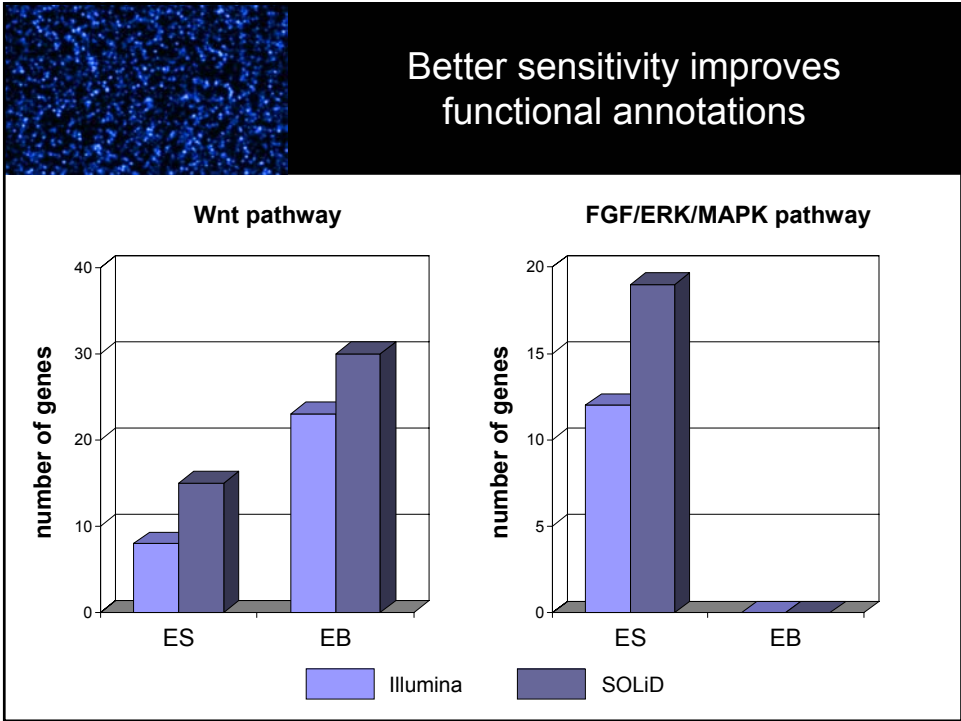
1. Random fragmentation of ribosomal depleted RNA
2. 1st strand cDNA synthesis using random hexamer-FDV primer
3. Addition of non-templated Cs by superscript II RTase
4. "Template switching" with RDV-riboGGG primer
5. 2nd strand cDNA synthesis & amplification
5. Emulsion PCR
6. SOLiD sequencing of ~150,000,000 beads per slide

Starting material: Ribo-depleted polyA and polysome associated RNA from mouse or human ES cells & EBs.

- Sequence ~100,000,000 mappable reads per library using SOLiD.
- Reads are 25, 30 or 35bp in length
- Make 3 libraries per biological state
- All tags are mapped to the genome first and then to a database of all known exon junctions.



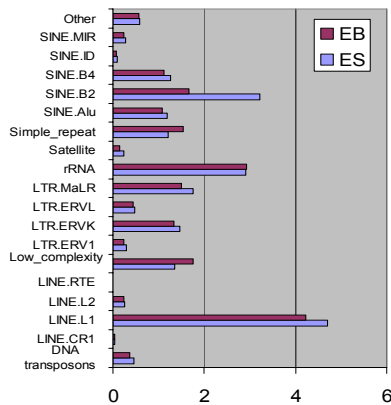
Up-regulated in EB library					Up-regulated in ESlibrary				
Refseq UCSC ID	Bstat	Fold change	Common name	Refseq ID	Refseq UCSC ID	Bstat	Fold change	Gene symbol	Refseq ID
16427	11.24	245.51	Hbb-bh1	NM_008219	577	10.26	21.43	Tcfcp2l1	NM_023755
2964	9.91	16.76	Car4	NM_007607	5456	9.77	41.83	Tdh	NM_021480
14226	9.47	43.17	Asb4	NM_023048	10018	9.54	13.69	Gm1967	NM_001033452
21211	9.11	28.00	Hand2	NM_010402	15887	8.95	12.79	Eg435970	NM_001034893
8843	8.77	20.70	Dkk1	NM_010051	985	8.33	19.32	Ifi202b	NM_008327
6847	8.53	7.94	Sl3ga6	NM_018784	10644	8.18	7.66	Mybl2	NM_008652
11406	8.40	42.28	Car14	NM_011797	6807	8.07	22.29	Dppa2	NM_028615
19444	7.51	13.57	Hmgb3	NM_008253	21838	8.01	15.36	Fgf4	NM_010202
18996	7.42	7.22	Pthr1	NM_011199	319	7.92	8.96	Mreg	NM_001005423
10966	7.35	7.65	Cldn11	NM_008770	15897	7.83	7.27	4933405K	NM_028913
2620	6.81	6.41	Tmem88	NM_025915	20848	7.82	20.54	Krt42	NM_0212483
2957	6.69	22.21	Lhx1	NM_008498	726	7.46	12.93	Nr5a2	NM_030676
16428	6.62	8.00	Hbb-y	NM_008221	5367	7.44	10.89	Tgm1	NM_019984
2677	6.60	7.64	Alox15	NM_009660	2133	7.42	14.45	Cobl	NM_172496
11726	6.55	5.51	Slc39a8	NM_026228	12114	7.38	23.83	Klf4	NM_010637
21817	6.50	10.40	Evx1	NM_007966	14257	7.34	7.26	Cav1	NM_007616
10049	6.46	9.75	Lmo2	NM_008505	11017	7.26	8.65	Phf17	NM_172303
18640	6.40	5.32	Dapk2	NM_010019	1041	7.23	7.13	Lefty2	NM_177099
2140	6.36	9.51	1500041B	NM_029861	1597	7.22	6.78	6330514A	NM_183152
19723	6.30	5.86	Nxf7	NM_130888	13160	7.21	10.56	Cnpy1	NM_175651
13035	6.30	5.88	C1qdc2	NM_026125	12236	7.11	6.54	Bnc2	NM_172870
11132	6.23	20.82	Tde2	NM_019911	1043	7.09	23.63	Lefty1	NM_010094
12231	6.20	43.36	Cer1	NM_009887	15996	7.08	12.92	Hsd17b14	NM_025330
3097	6.15	5.24	Skap1	NM_001033186	7789	7.00	12.23	1700061G	NM_030141
6586	6.14	6.38	Klhl6	NM_183390	21186	6.86	8.41	E130014J	NM_001040400
3257	6.13	7.34	Ramp2	NM_019444	17531	6.85	10.96	Cign	NM_008904
11243	5.83	6.17	Etna1	NM_010107	16769	6.80	29.50	Mylpf	NM_016754
7976	5.83	4.00	Ak220484	NM_001083628	7565	6.79	7.14	Zfp57	NM_001013745



Transcriptome Discovery: in ES cells and EBs

Category	ES	EB
Known exons	74.2%	75.5%
Predicted exons	5.6%	7.8%
Known regions	7.9%	5.6%
Predicted regions	4.8%	3.9%
Conserved regions	2.8%	3.2%
Other regions	5.6%	7.8%

Widespread and dynamic Repeat expression:



>300 Repeat elements display bidirectional expression in the ES cell state

Having established Ecol as a critical factor for the generation of DSB-induced cohesion, we asked whether the failure to generate cohesion in aged G_2M cells results from limiting activity. To test this hypothesis, we treated Ecol during G_2M in the absence of shadD DNAs. Indeed, overproduction of but not *ecol*⁺ bypasses the requirement for shadD to generate cohesion in G_2M (Fig. 3C). Ecol and *ecol*⁺ are present at similar (Fig. 3D) and *ecol*⁺ rescues a function in because it complements the *ecol*⁻303 operative temperature (Fig. 3D). These results suggest that in G_2M , Ecol acetylation activity is limiting in undamaged cells and, given its DNAs, this activity is elevated in the DNA damage checkpoint.

As we show that the generation of cohesion M is Ecol dependent but replication is not (also reported in 145). This contradicts our model, which posits that cohesion proteins can only occur in the context of DNA replication, and Ecol (C) directly allows the replication fork through the cohesion ring in 3-phase after, we suggest that Ecol directly converts gamma-hemolysin-cohesion complex to its active state. During S phase, Ecol associates with components 1', 1'', and the allows to establish cohesion before the onset of

Developmentally Regulated Activation of a SINE B2 Repeat as a Domain Boundary in Organogenesis

Victoria V. Lopez,^{1,2} Gratian G. Pfeifer,^{1,2} Susumu Hahn,^{1,2} Florian Cramer,^{1,2} Bong-Gook Ju,¹ Kenneth A. Duhl,¹ Kary Hart,¹ Eusebio Ray,¹ Angel Garcia-Diaz,¹ Xiangyan Zhu,¹ Yan Yang,¹ Liuhua Mastrolia,¹ Christopher K. Glass,¹ Michael G. Rosenfeld^{1*}

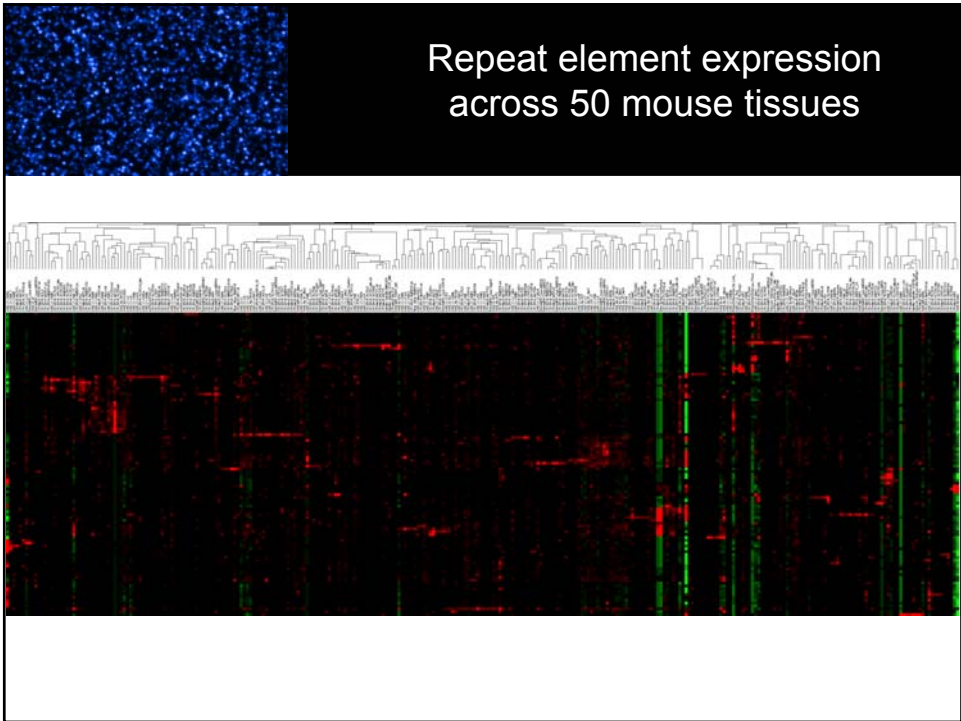
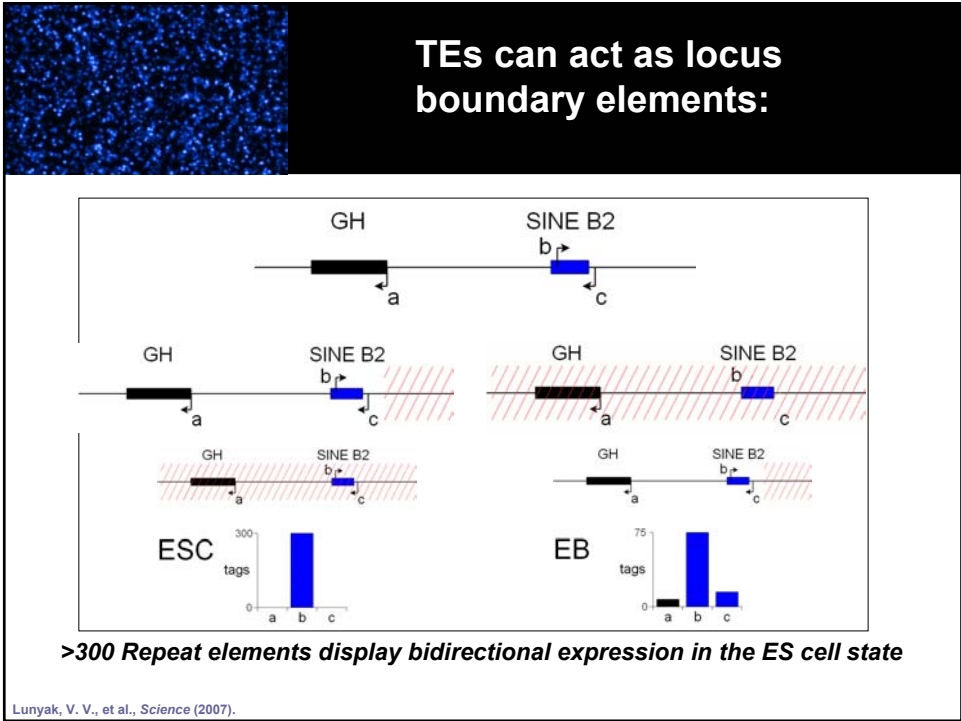
The temporal and spatial regulation of gene expression in mammalian development is linked to the establishment of functional chromatin domains. Here, we report that tissue-specific transcription of a retrotransposon repeat in the murine growth hormone locus is required for gene activation. This repeat serves as a boundary to block the influence of repressive chromatin modifications. The repeat element is able to generate short, overlapping Ptd II- and Ptd III-driven transcripts, both of which are necessary and sufficient to enable a restructuring of the regulated locus into nuclear compartments. These data suggest that transcription of interspersed repetitive sequences may represent a developmental strategy for the establishment of functionally distinct domains within the mammalian genome to control gene activation.

The growth hormone (*GH*) gene provides a well-studied transcription unit that is highly regulated for defining low-specific chromatin modifications (*1-10*) might be responsible for the spatial and temporal order of lineage specification events in the developing pituitary gland. The human *GH* locus is represented by a cluster of five *GH*-related genes that are regulated by a Ptd I-

13 JULY 2007 VOL 317 SCIENCE www.sciencemag.org



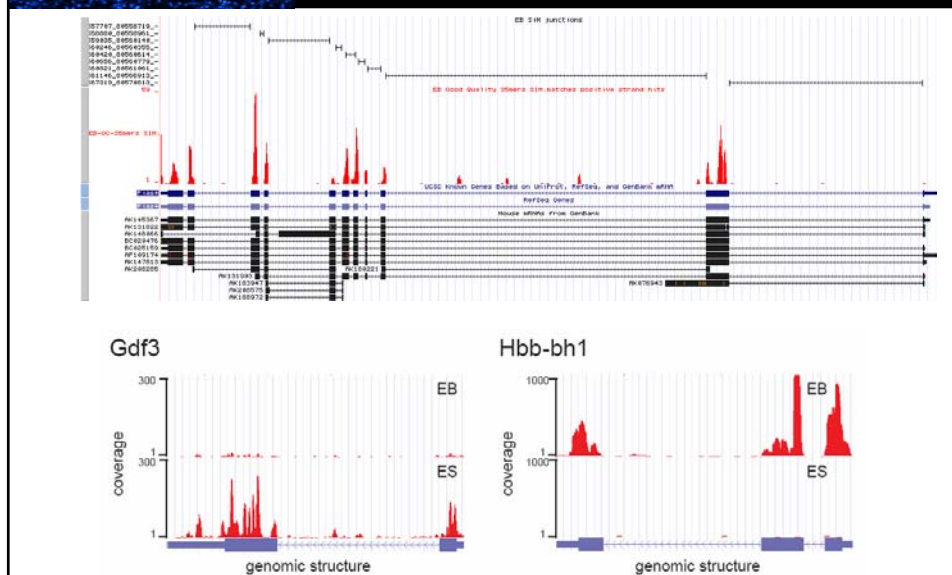
>300 Repeat elements display bidirectional expression in the ES cell state



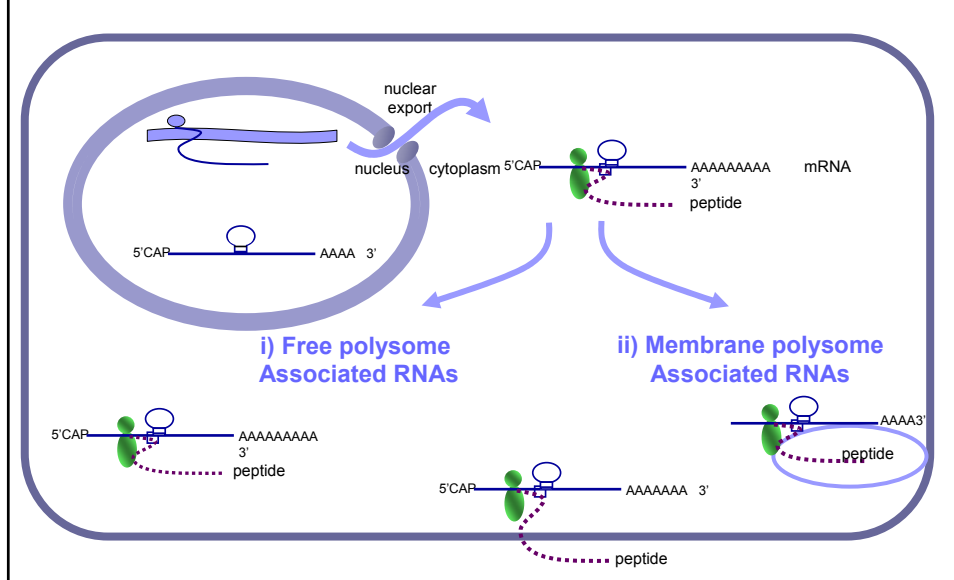
TE can act as alternate promoters

- ❖ TEs known to provide alternative promoters to nearby genes (Speek et al, Genomics, 2008).
- ❖ CAGE delineated >100,000 TE promoters
 - <100kb upstream of a RefSeq
 - same strand as a RefSeq
 - not overlapping a RefSeq
- ❖ >700 confirmed as alt. promoters by ESTs
- ❖ Confirmation of novel promoters RT-PCR, cDNA sequencing and RACE

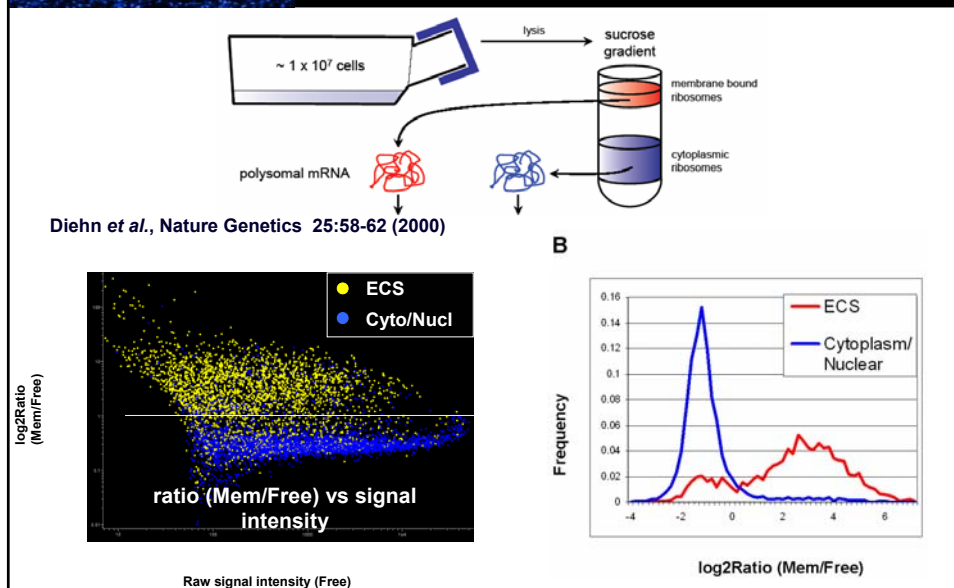
Viewing exon and junction expression in a genomic context:



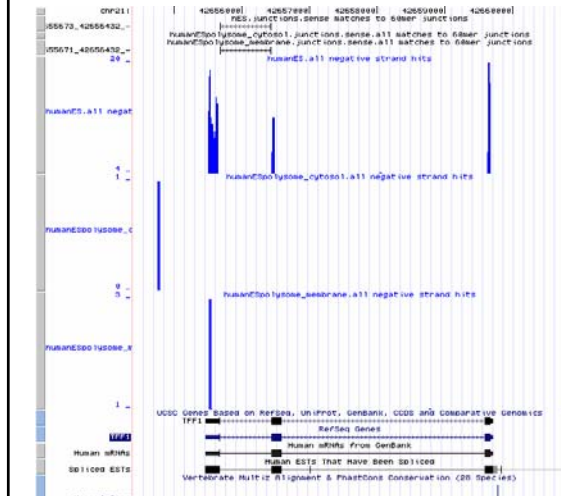
The Polysome-associated transcriptome:



hES PolyA Vs Polysome RNA sequencing:

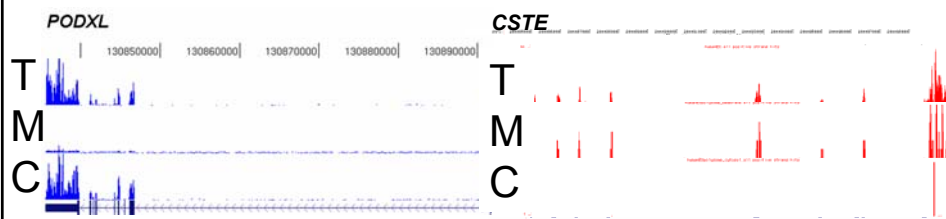


Expression Vs translation



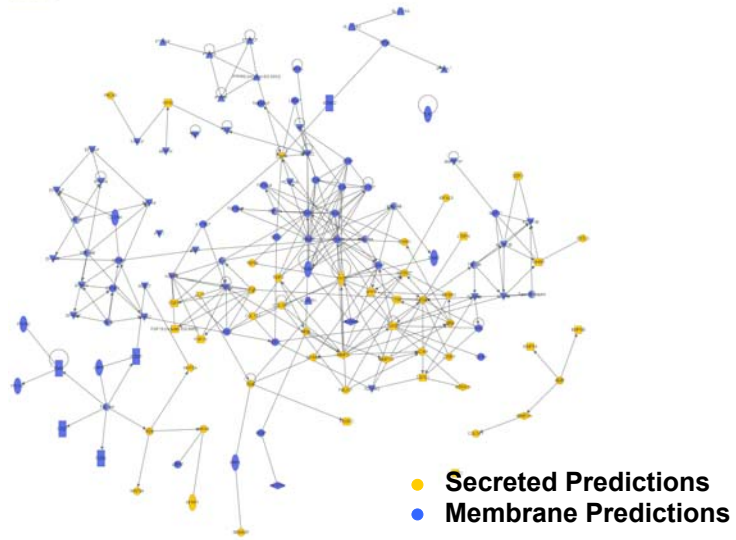
Tag count
 Total: 1115
 Cytosol: 0
 Membrane: 3

hES PolyA Vs Polysome RNA sequencing:

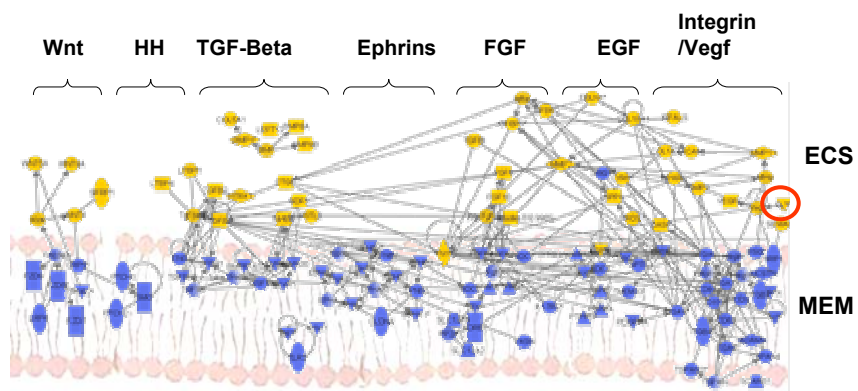


	Total	Cytosol	Membrane
Expressed Refseq genes	12,497	9,387	6,700
All aceview diagnostic variants	53,190	26,005	11,060
All aceview diagnostic loci	30,281	17,851	9,160
Ratio (variants/loci)	1.76	1.46	1.21

Network analysis of ECS signalling:



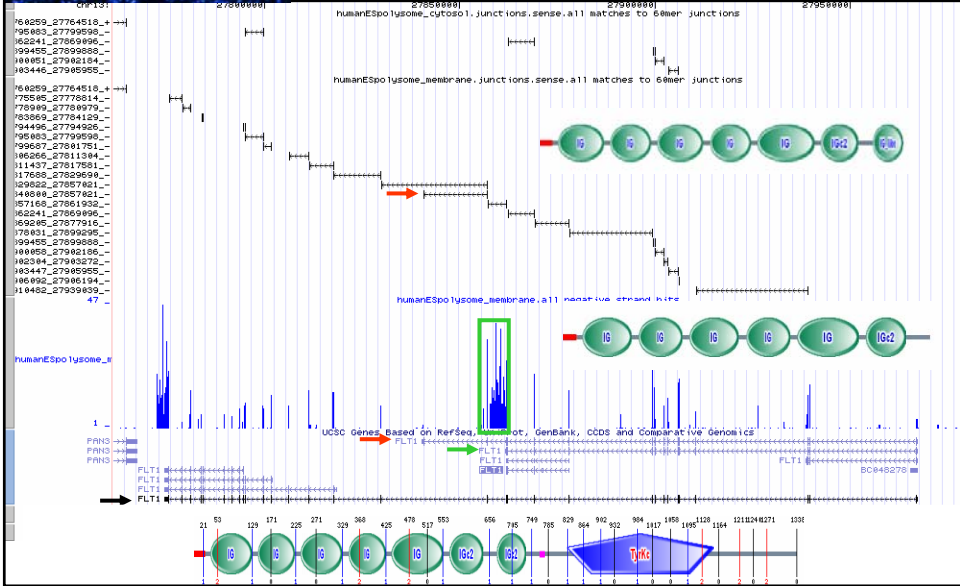
Modelling the extracellular space in ES cells:



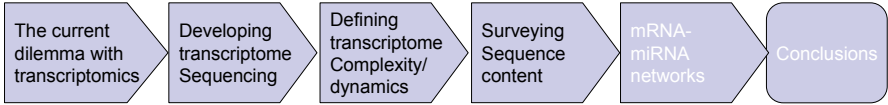
Ingenuity Pathways Analysis™

- Secreted Predictions
- Membrane Predictions

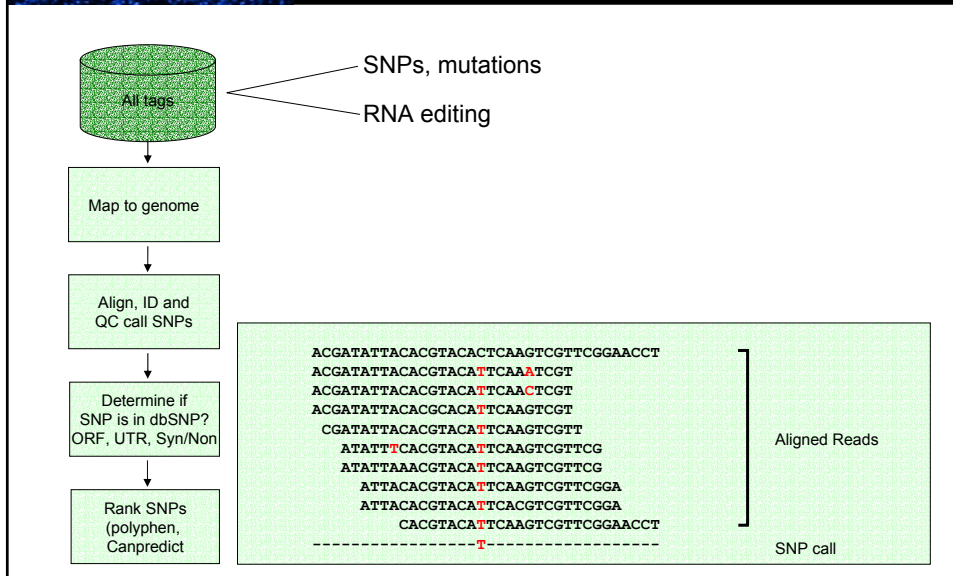
Transcription output of VEGFR1 In hES:



Presentation overview



Screening sequences for substitutions (expressed SNPs, mutations, RNA editing)



Validation of sequence variations

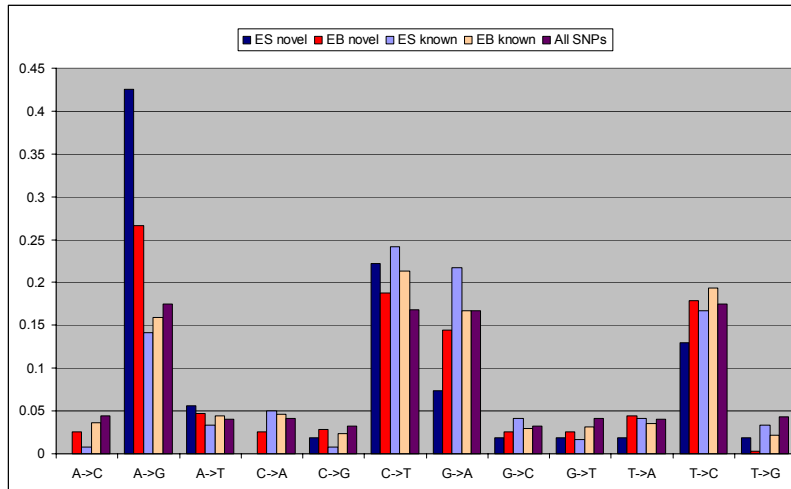
Chromosome	Position	Strand	Gene	Predicted SNP	Amino acid change	Validated
chr1	174437452	+	Tagln2	C->T	UTR	Yes
chr1	34501748	+	Imp4	C->T	UTR	Yes
chr11	84777618	+	Car4	C->T	A->V	Yes
chr14	22616785	+	Samd8	A->G	UTR	Yes
chr18	38460154	+	Rnf14	C->T	A->V	Yes
chr19	41992404	+	Pgam1	C->T	UTR	Yes
chr2	119453391	+	Nusap1	A->G	K->R	Yes
chr3	121978860	+	Dnttip2	T->A	L->Q	Yes
chr3	93330080	+	St00a1.1	C->T	UTR	Yes
chr4	15903284	+	Nbn	C->G	Y->F	Yes
chr5	138548192	+	Zkscan1	A->G	UTR	Yes
chr5	138548273	+	Zkscan1	A->G	UTR	Yes
chr7	38752418	-	C80913	T->C	N->D	Yes
chr7	50209670	-	EG668668	C->G	G->R	Yes
chr7	5076892	-	Rasl2-9	T->A	E->D	No
chr9	123371057	+	Lars2	T->G	UTR	Yes
chr9	57681529	-	Arid3b	G->T	Q->K	Yes
chr9	64083463	+	Uchl4	A->T	D->E	No

Dnttip2; T->A; L269Q

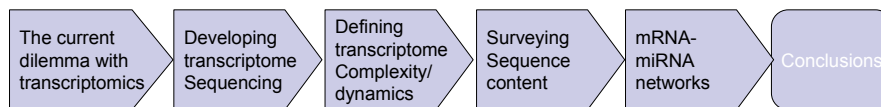
Arid3b; G->T; Q->K

Uchl4; T->A; Y->F

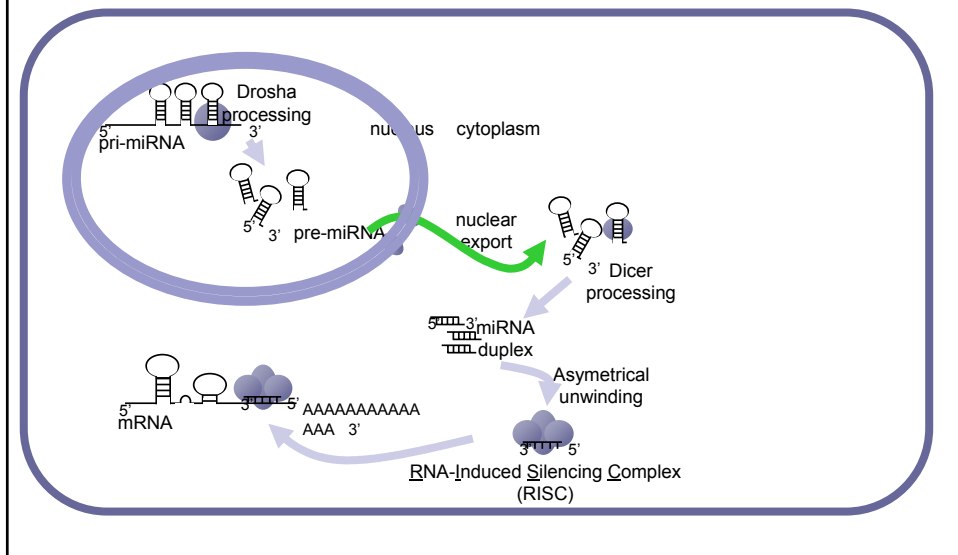
SNP detection in ES: Evidence of RNA editing?



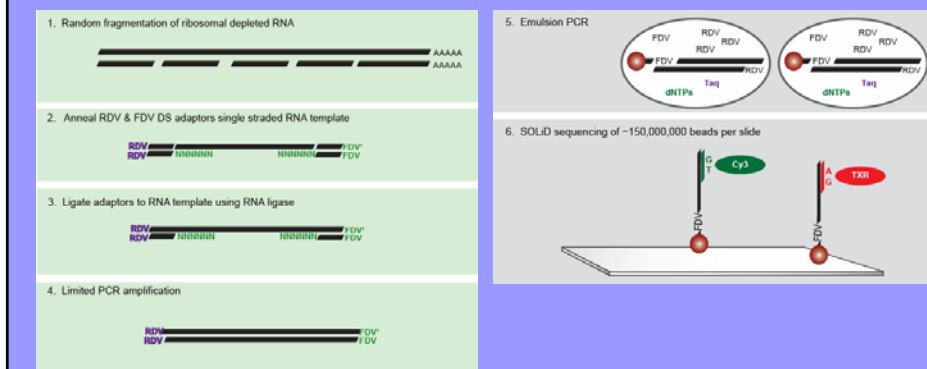
Presentation overview



miRNA suppression of mRNA targets

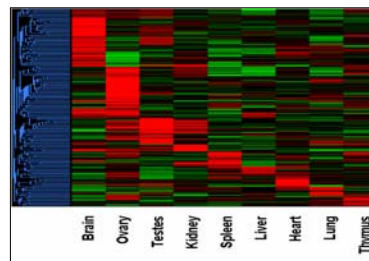
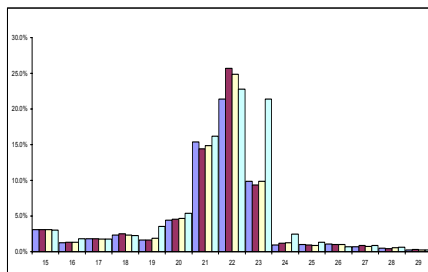
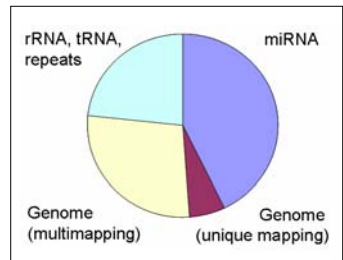
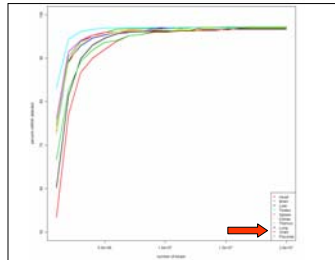


miRNA sequencing in mammalian ESCs & EBs



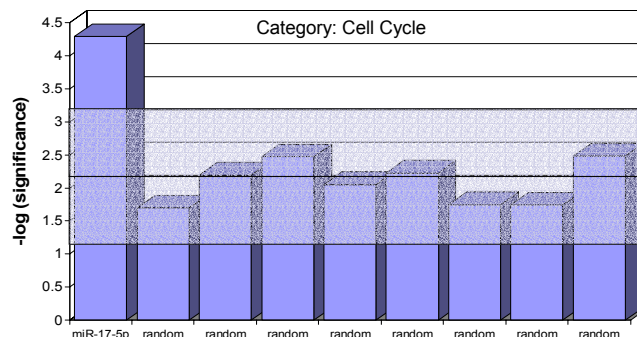
Sequence ~10,000,000 15-35bp reads using SOLiD
 Map against miRBase, then the genome
 Cluster miRbase matches and summarise isomiR complexity

miRNA sequencing in mammalian ESCs & EBs

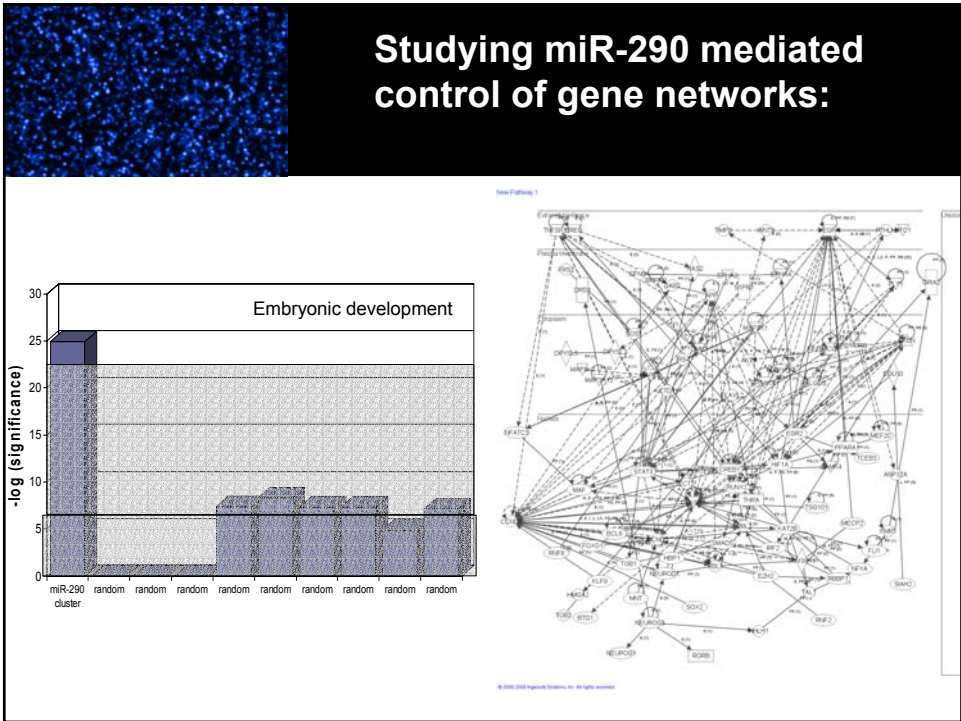
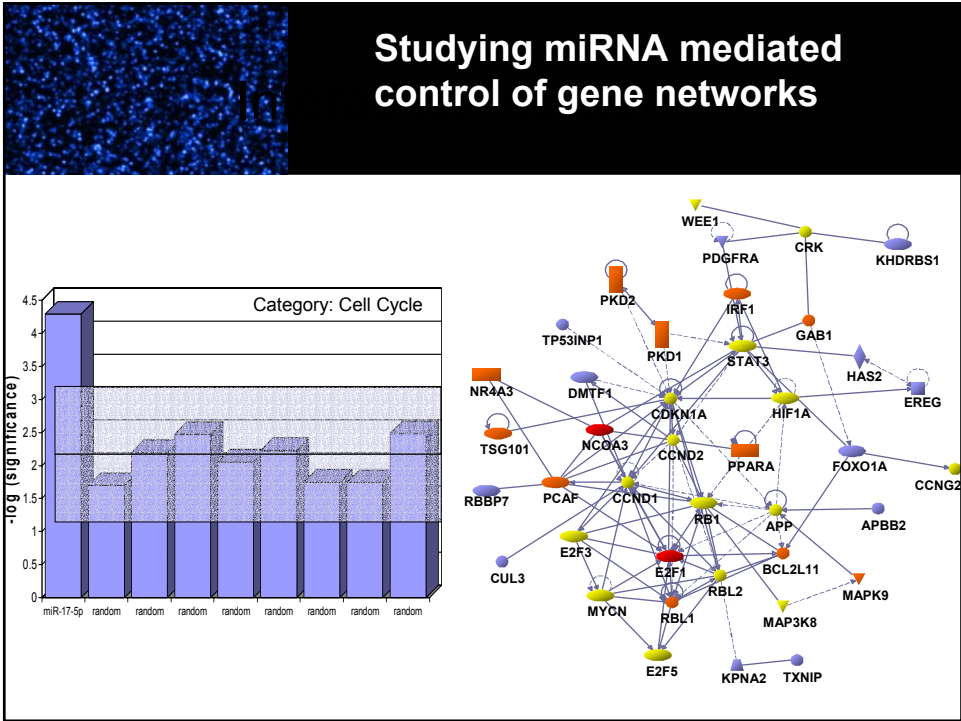


Studying miRNA mediated control of gene networks:

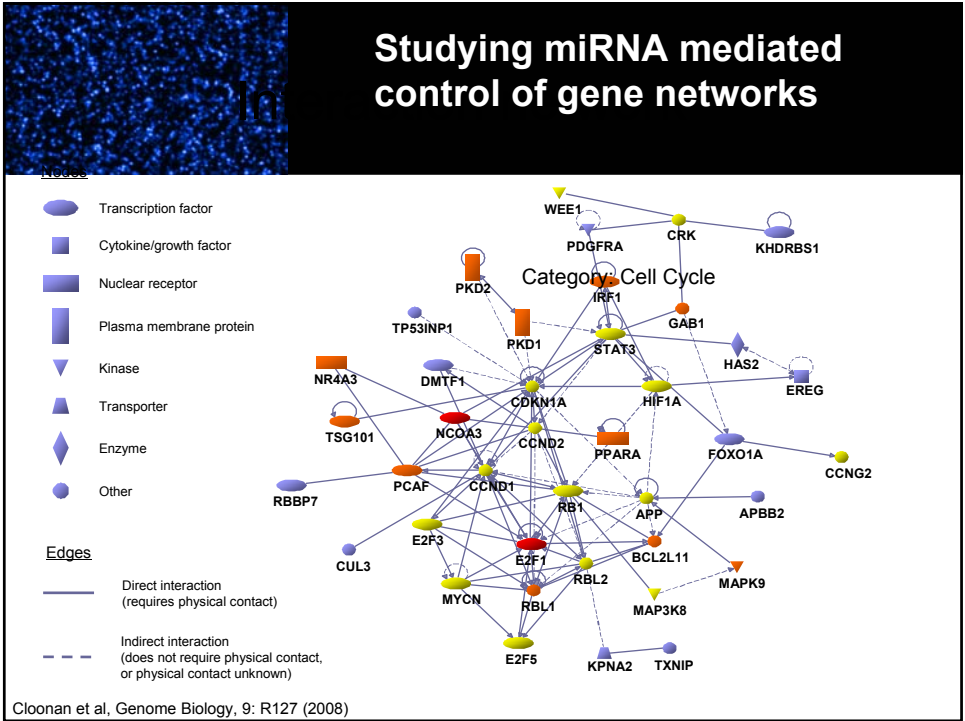
1. Define the miR targets of the polycistronic miRs (PicTar)
2. Determine functional relationship of targets (IPA)
3. Benchmark against similarly size random sets of targets



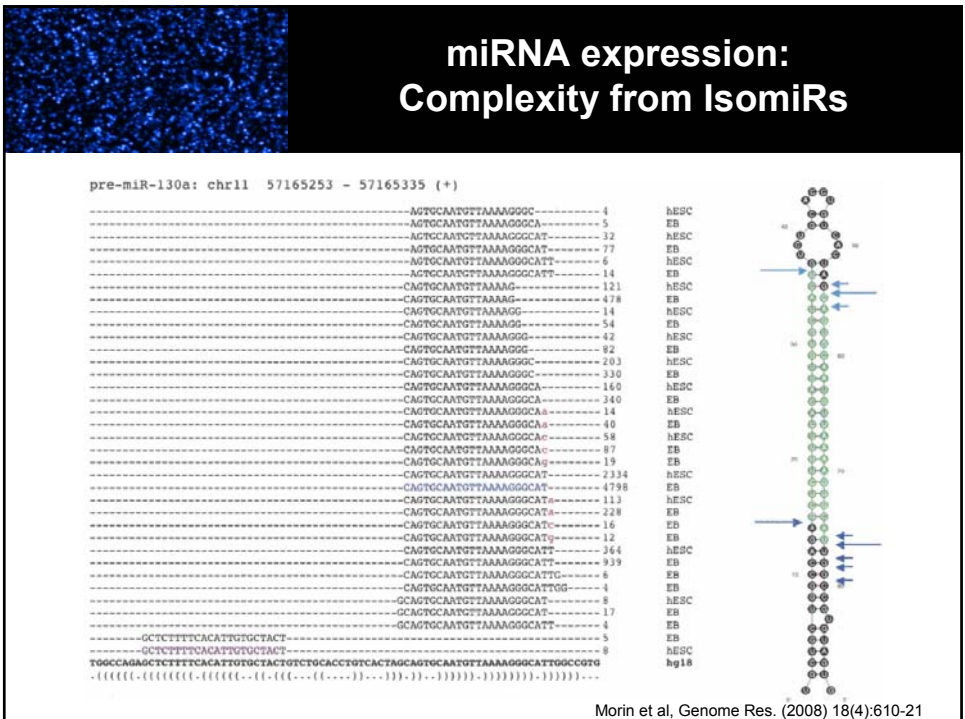
Cloonan et al, Genome Biology, 9: R127 (2008)



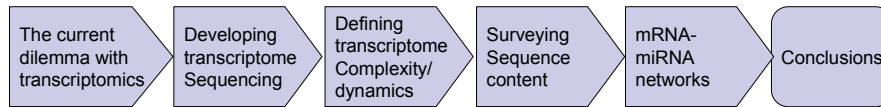
Studying miRNA mediated control of gene networks



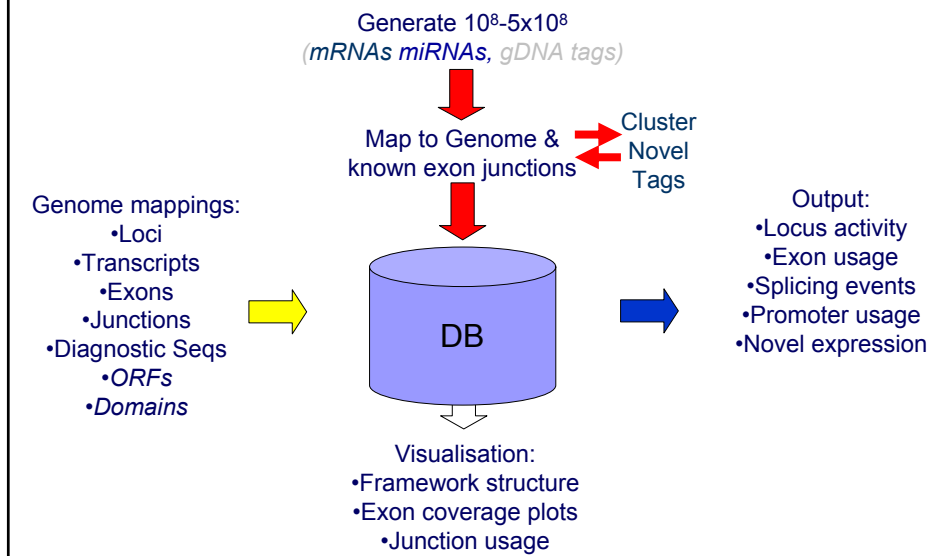
miRNA expression: Complexity from IsomiRs



Presentation overview



Summary:



Conclusions:

RNAseq is a powerful tool for studying RNA abundance, dynamics and complexity with unprecedented sensitivity

Transcriptome discovery is identifying roles for novel expression

Polysome associated RNAseq gives novel insights into those transcripts that are actively translated.

SQRL provides a complete view of small RNA complexity and abundance.

It is also now possible to screen the transcriptome for expressed SNPs, mutations (mis-sense, and indels), allelic specific expression and RNA editing.

Challenges:

-Sequence content (more reads, longer reads, more accurate sequencing).

-Ambiguity of mapping. "Black holes" in the transcriptome.

-Transcript length affect shotgun tag sensitivity.

-Matching to a reference is easy, de novo assembly is hard.

-Diagnostic sequences only infer transcription complexity

-Novel combinations of junction usage (ie alt splicing event between exon 1 & 3 cannot be conclusively tied to alternate 3'utr usage.)

SOLiD:

Kevin McKernan
Clarence Lee
Heather Peckham
Stephen McLaughlin

SREK:

Scott Kuersten
Jian Gu
Catalin Barbacioru

IMB array facility:

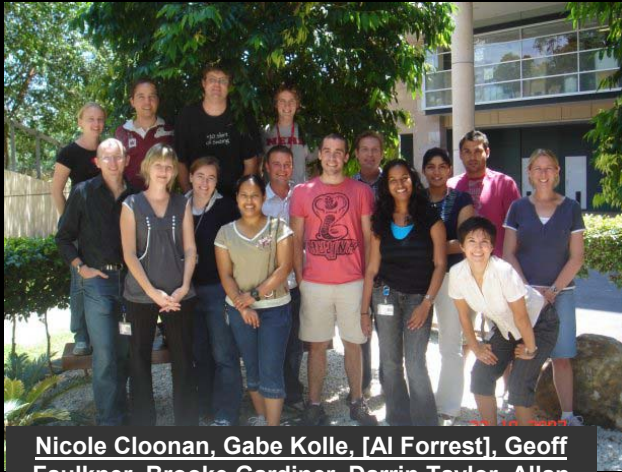
Tina Maguire
Keerthana Krishna

Stem Cell Research:

Andrew Perkins
Steve Bruce
Andrew Laslett
George Zhou

Repeat analysis:

Piero Carninci
Shintaro Katayama
Valerio Orlando



Nicole Cloonan, Gabe Kolle, [Al Forrest], Geoff Faulkner, Brooke Gardiner, Darrin Taylor, Allan Robertson, Shivangi Wani, Graham Bethel.

