

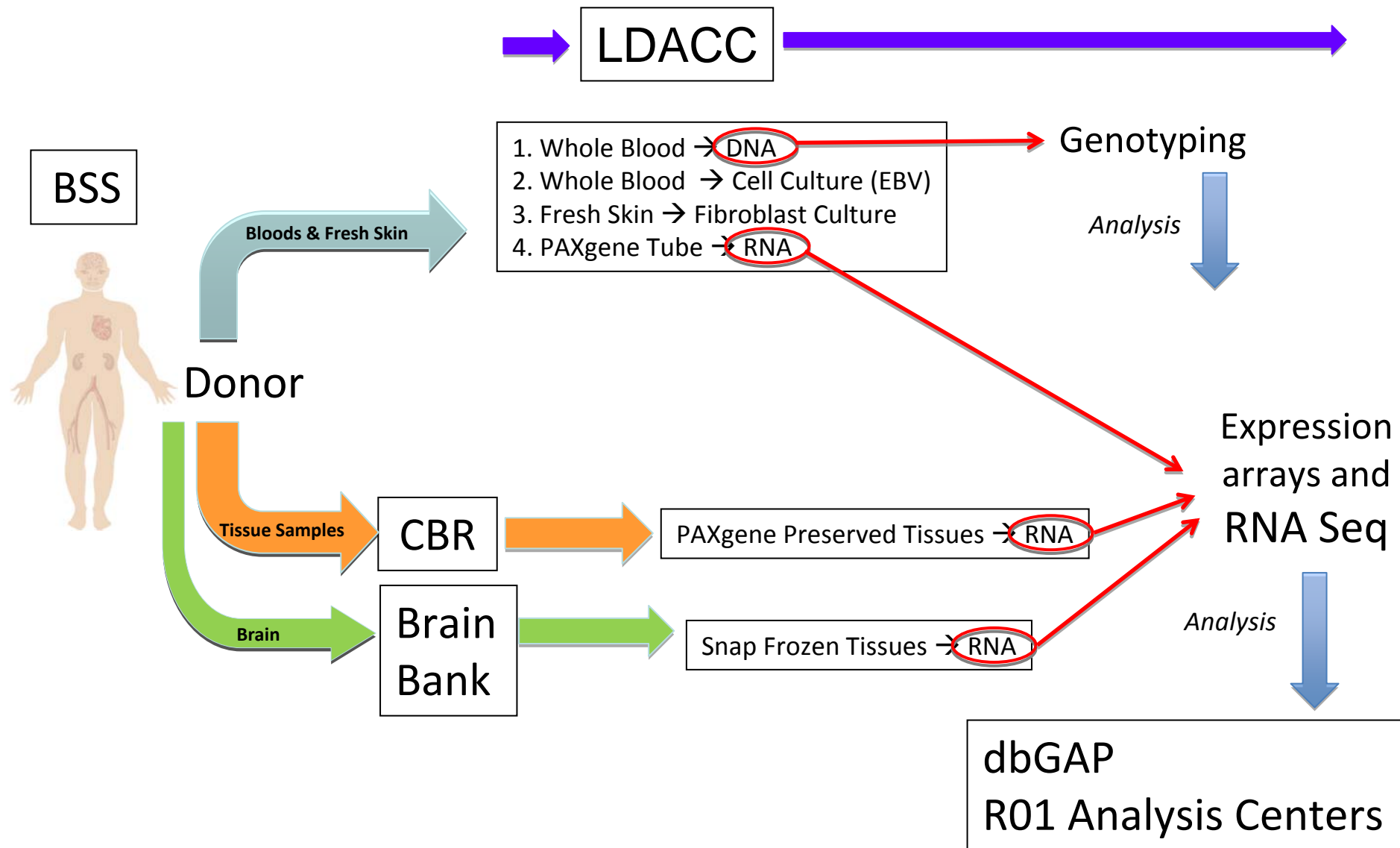


GTE_x – Sample Molecular Quality and Data Production

Kristin Ardlie

23 February, 2012

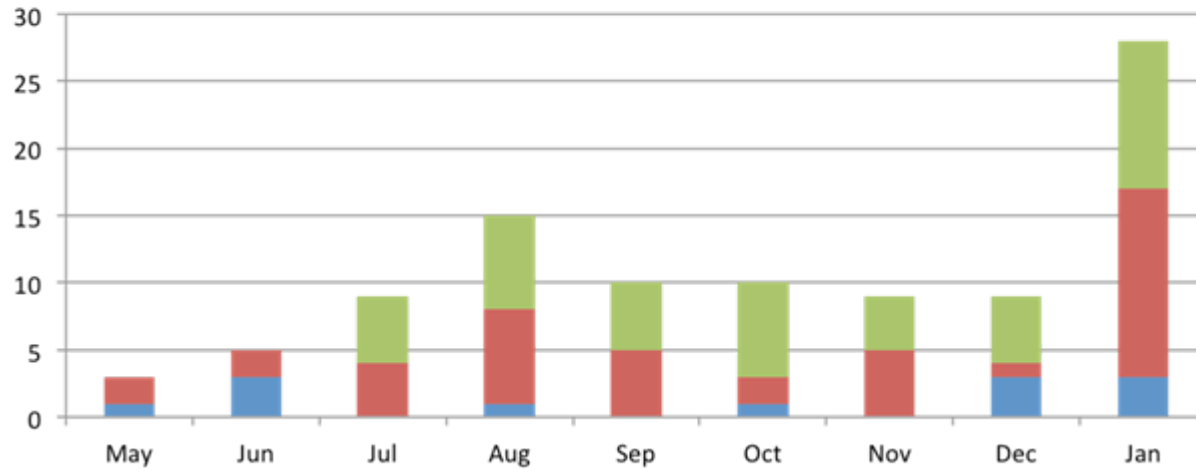
GTEx Sample Workflow



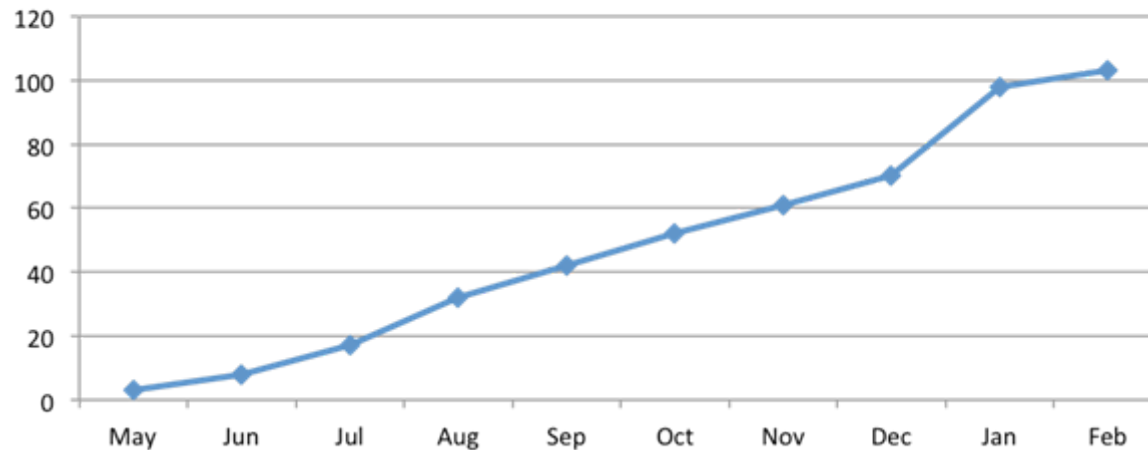
Sample Receipt by LDACC



Sample Receipt by Month

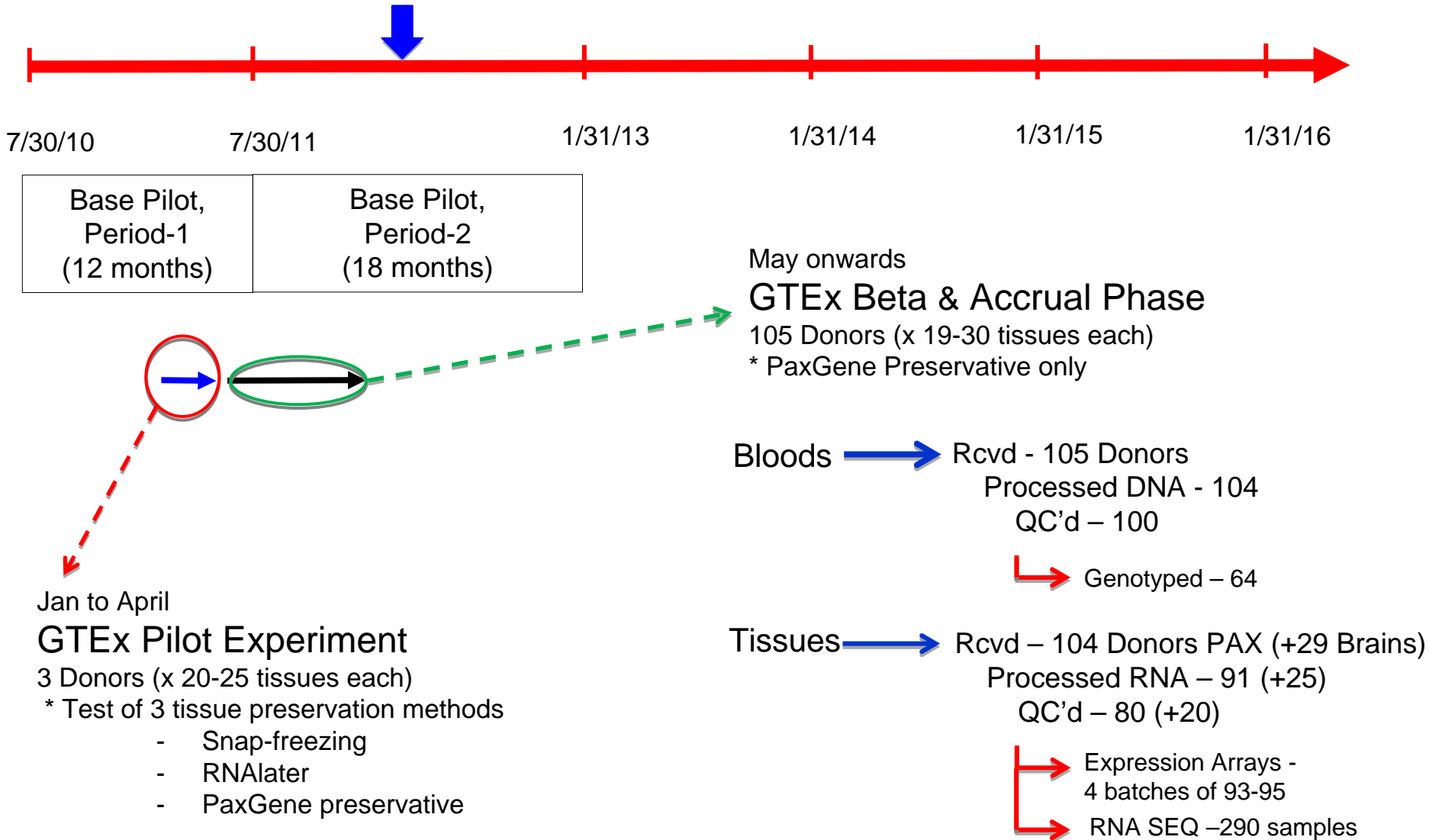


Sample Receipt Cumulative



N=105 donors

GTEx Project Current Status



Target Tissues

Paxgene Preserved Tissues

- Adipose
- Adrenal Gland
- Artery – Aorta
- Artery – Coronary
- Artery – Tibial
- Bladder
- Breast - Mammary Tissue
- Cervix – Ecto & Endocervix
- Colon
- Esophagus – Mucosa & Muscularis
- Fallopian Tube
- Heart
- Kidney – Cortex & Medulla
- Liver
- Lung
- Muscle – Skeletal
- Nerve – Tibial
- Ovary
- Pancreas
- Pituitary

- Prostate
- Skin
- Spleen
- Stomach
- Testis
- Thyroid
- Uterus
- Vagina

Fresh Frozen Brains

- Brain – Cerebellum
- Brain – Cortex
- Cerebellar Hemisphere
- Frontal Cortex (BA9)
- Hippocampus
- Substantia nigra
- Anterior cingulate cortex (BA24)
- Amygdala
- Caudate (basal ganglia)
- Nucleus accumbens (basal ganglia)
- Putamen (basal ganglia)
- Hypothalamus
- Spinal cord (cervical c-1)

Currently collecting 9-
30 tissues per donor

Target Tissues – High Priority



Paxgene Preserved Tissues

- **Adipose**
- Adrenal Gland
- Artery – Aorta
- Artery – Coronary
- **Artery – Tibial**
- Bladder
- Breast - Mammary Tissue
- Cervix – Ecto & Endocervix
- Colon
- Esophagus – Mucosa & Muscularis
- Fallopian Tube
- **Heart**
- Kidney – Cortex & Medulla
- Liver
- **Lung**
- **Muscle – Skeletal**
- **Nerve – Tibial**
- Ovary
- Pancreas
- Pituitary

- Prostate
- **Skin**
- Spleen
- Stomach
- Testis
- **Thyroid**
- Uterus
- Vagina

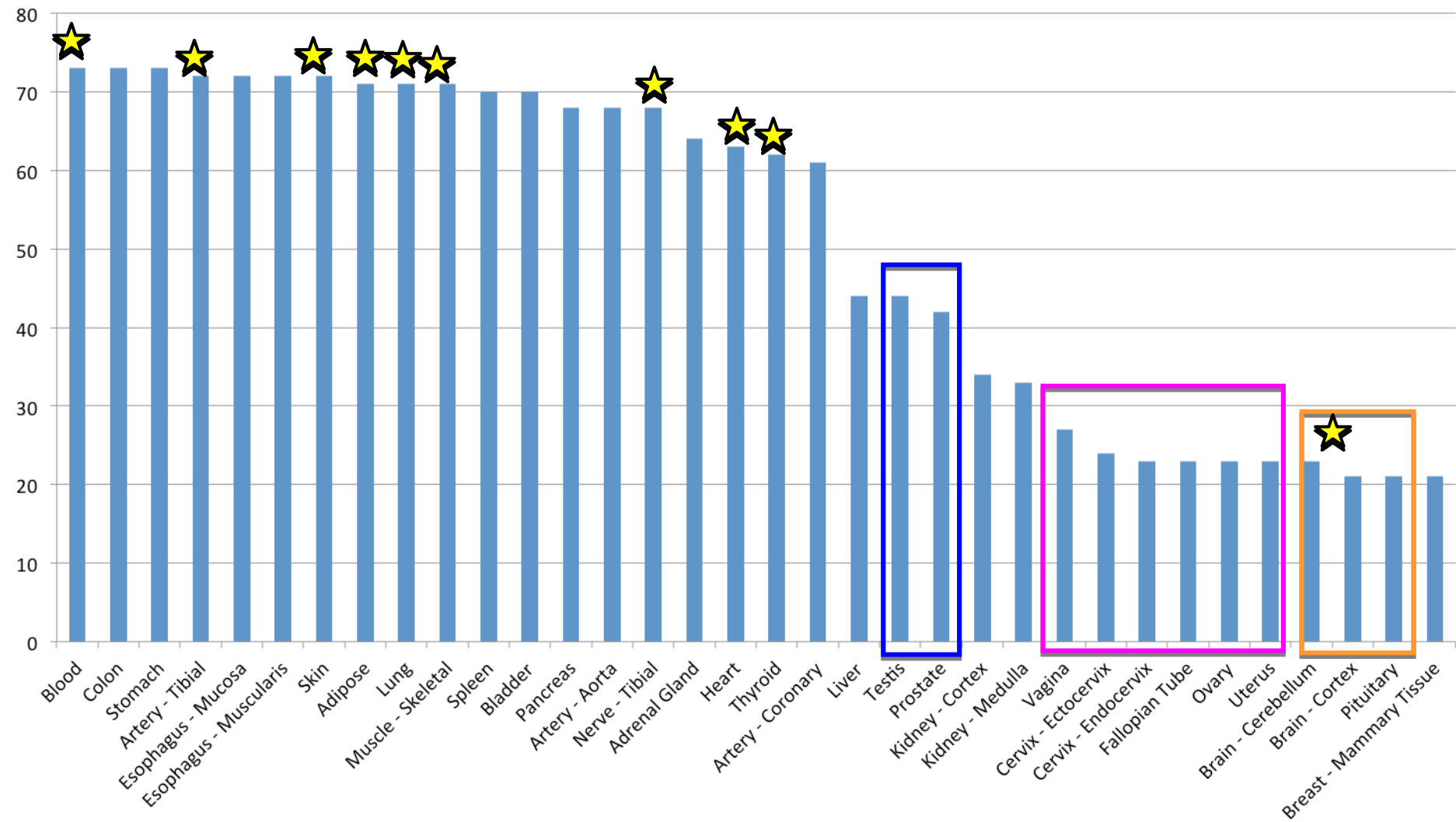
- **Fibroblast Cell Line**
- **Blood**
- **Lymphoblastoid Cell Line**

Fresh Frozen Brains

- **Brain – Cerebellum**
- **Brain – Cortex**
- **Cerebellar Hemisphere**
- **Frontal Cortex (BA9)**
- **Hippocampus**
- **Substantia nigra**
- **Anterior cingulate cortex (BA24)**
- **Amygdala**
- **Caudate (basal ganglia)**
- **Nucleus accumbens (basal ganglia)**
- **Putamen (basal ganglia)**
- **Hypothalamus**
- **Spinal cord (cervical c-1)**

21 Tissues + 3 = 24 Priority

Tissue Distribution

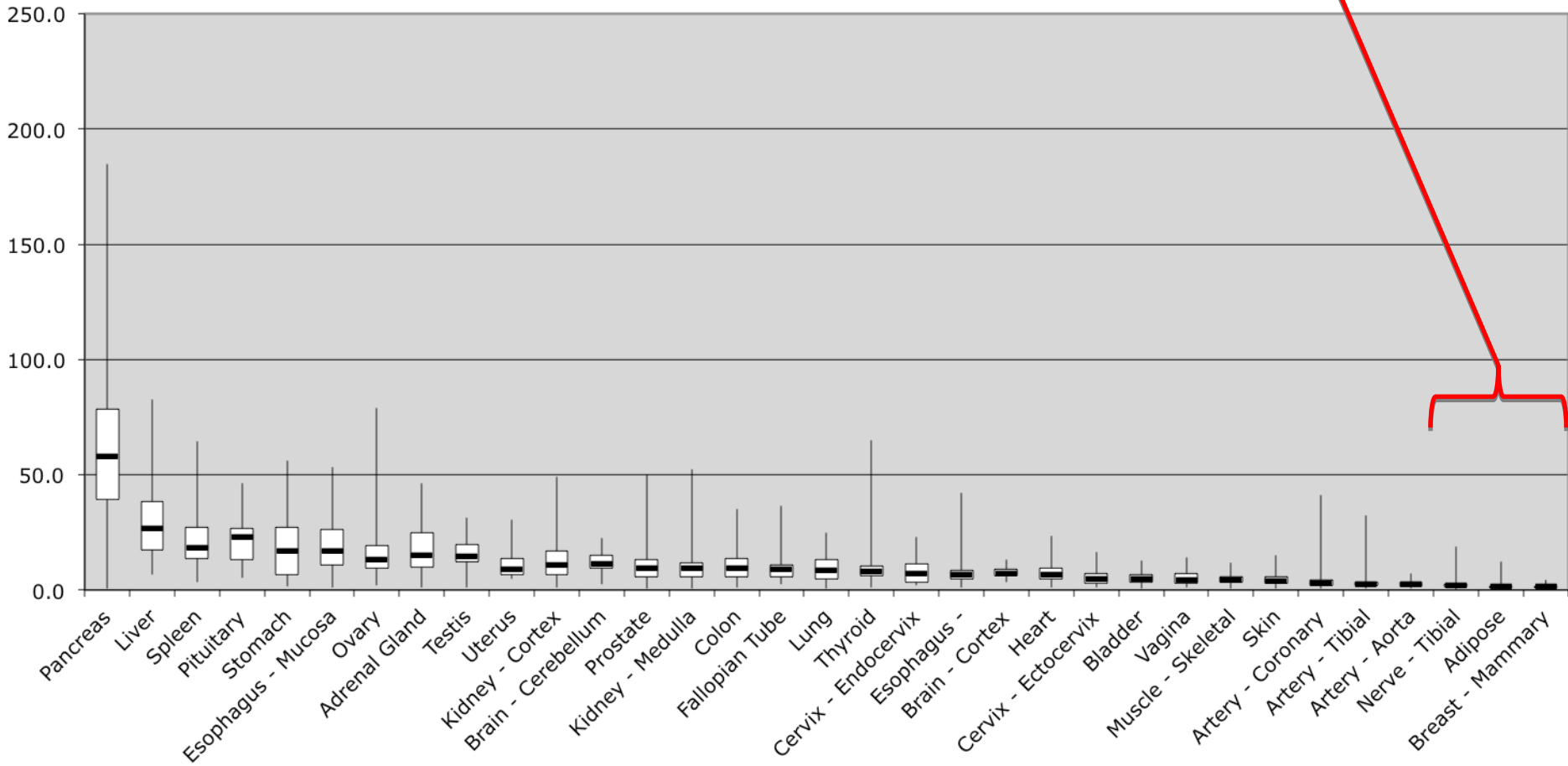


N = 75 Donors

RNA Yields

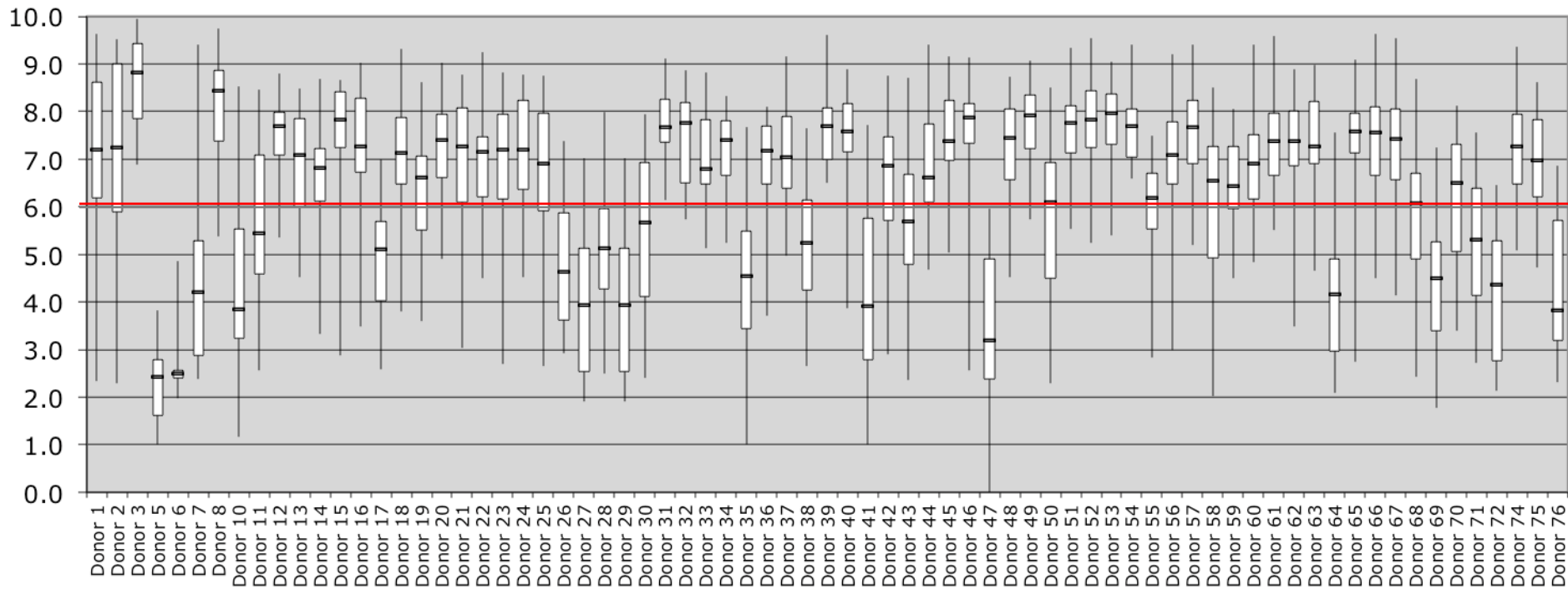


1-3ug



N = 1667 Paxgene Tissues

RNA Quality – Paxgene Tissues



N=75 donors

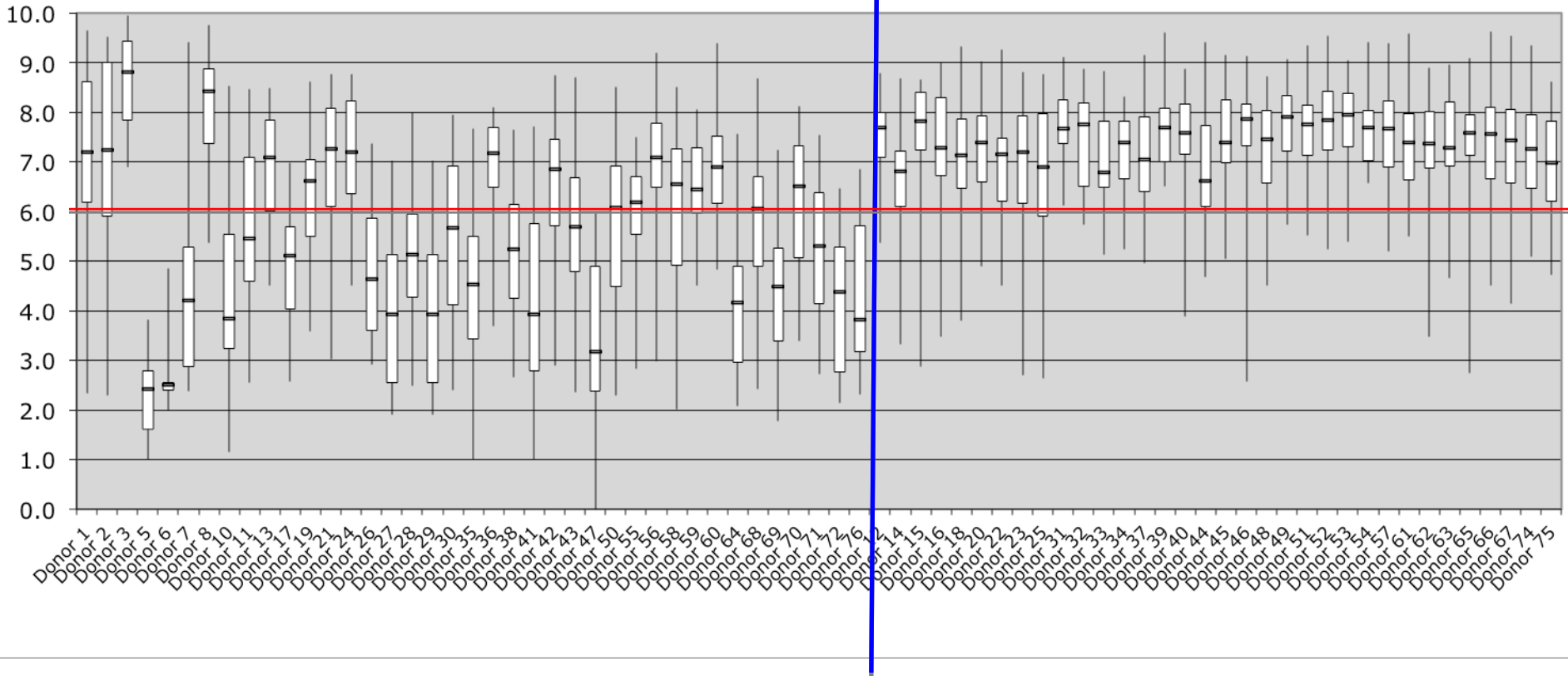
RNA Quality varies by Donor

RNA Quality – Paxgene Tissues



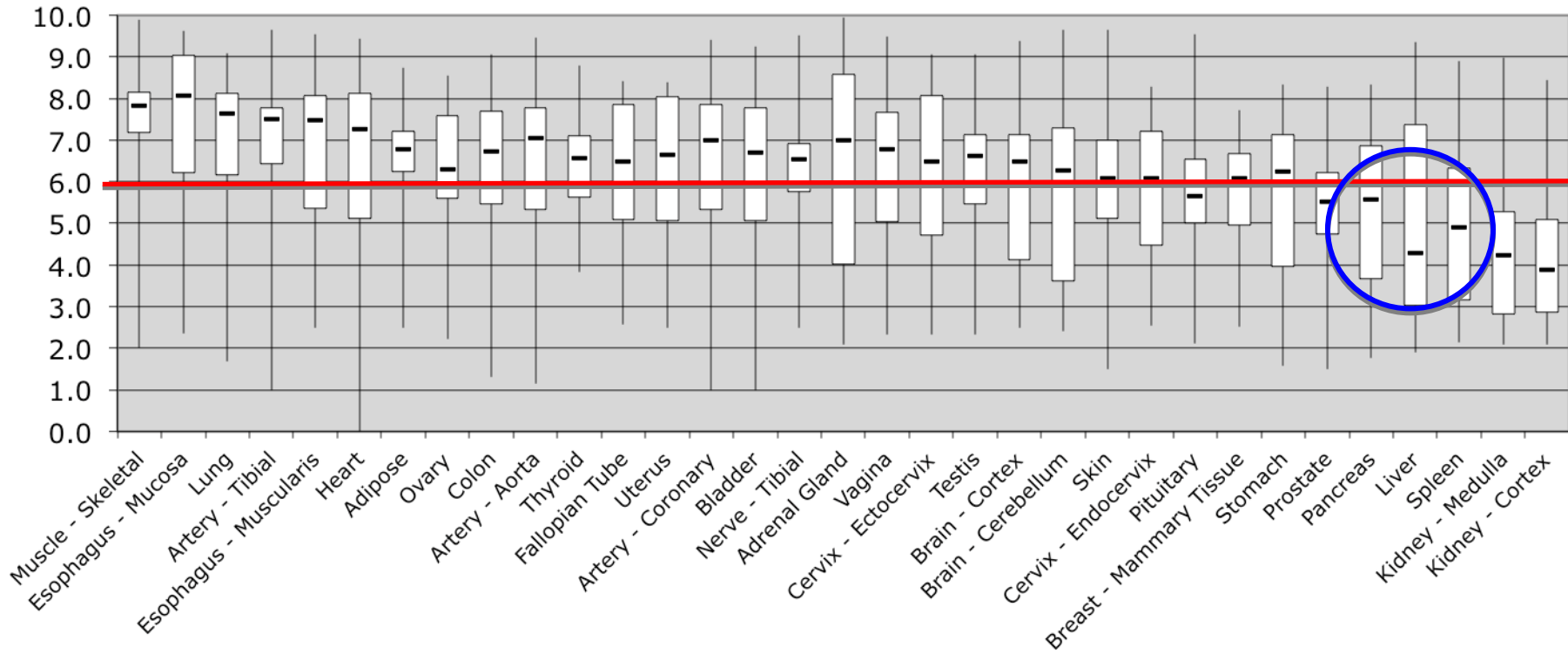
N = 41 POSTM Donors

N = 34 OPO Donors



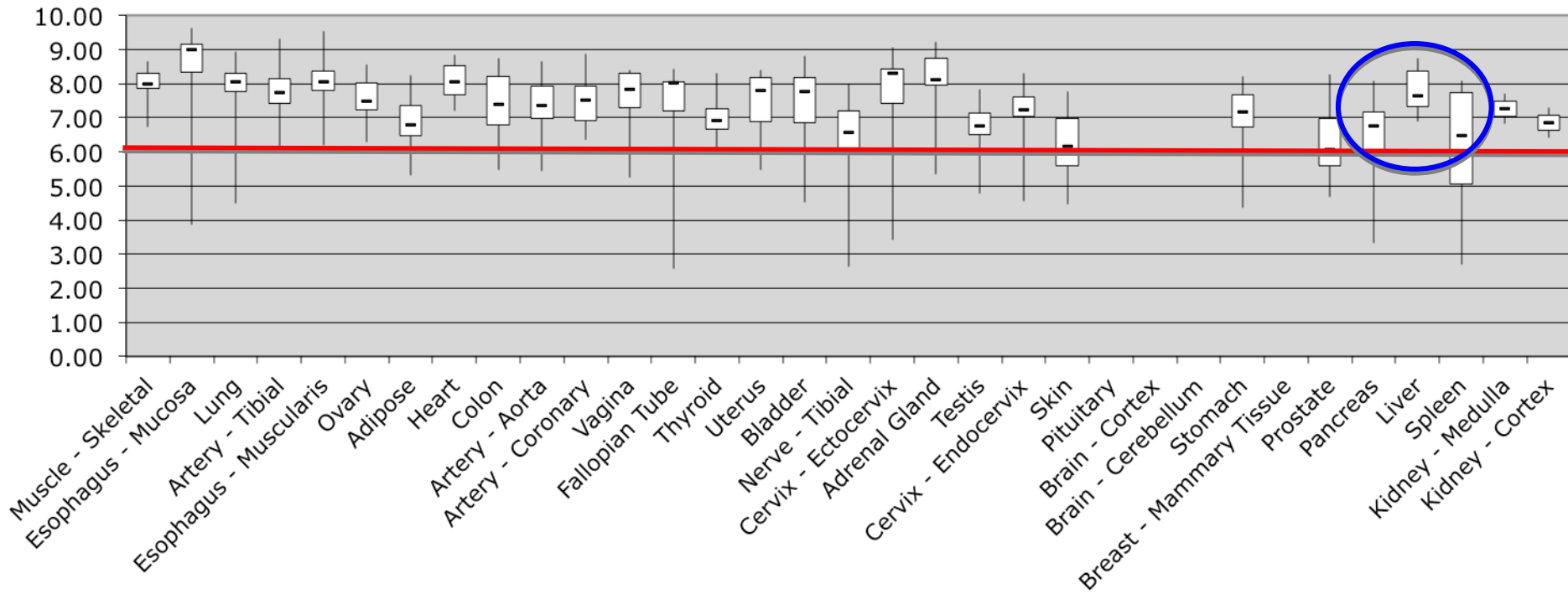
Donor RNA variation influenced by collection type

RNA Quality by Tissue Type



N=1667 tissues

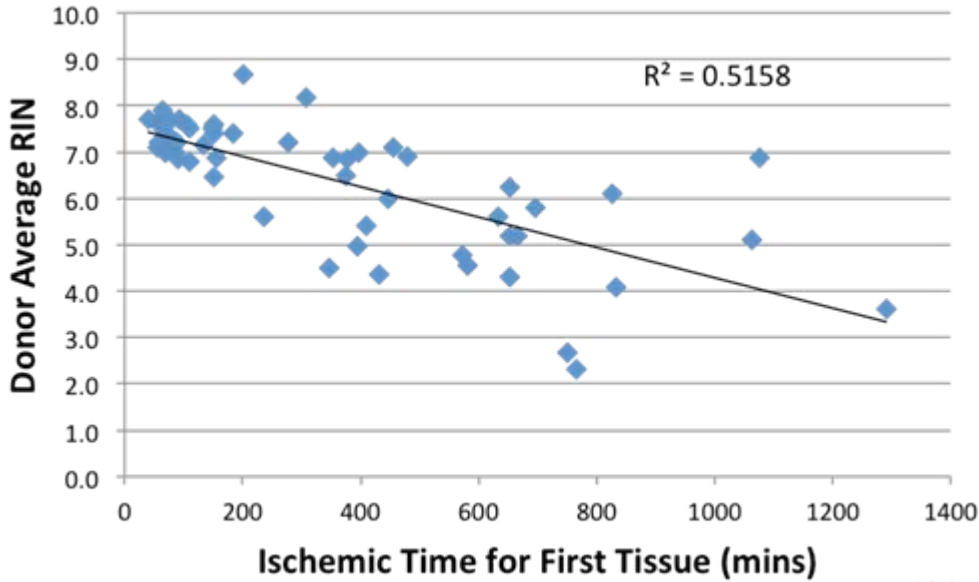
RNA Quality by Tissue Type



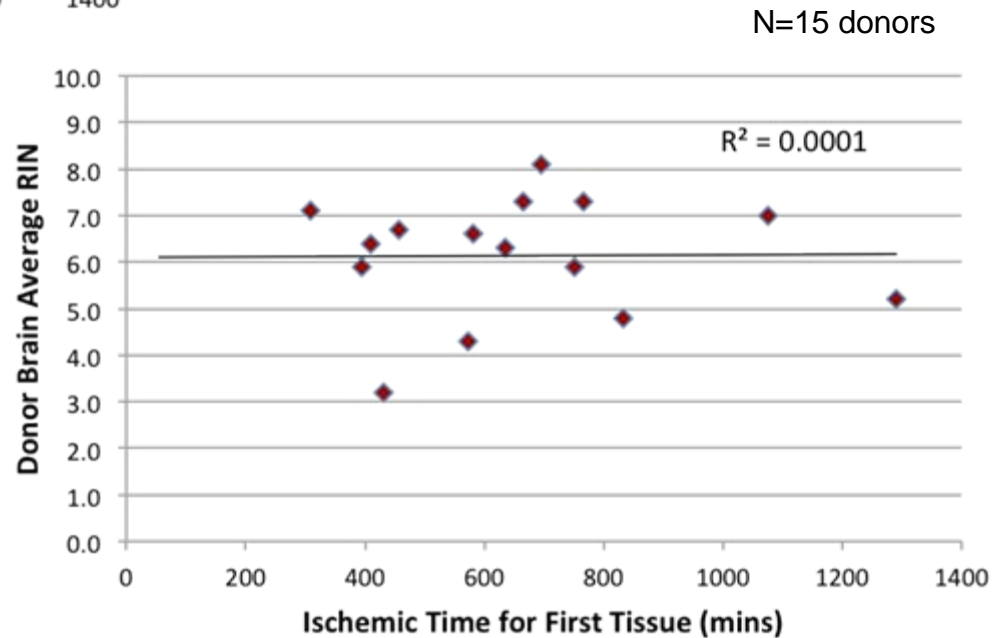
OPO donors only

Tissue-specific RNA quality also donor-driven

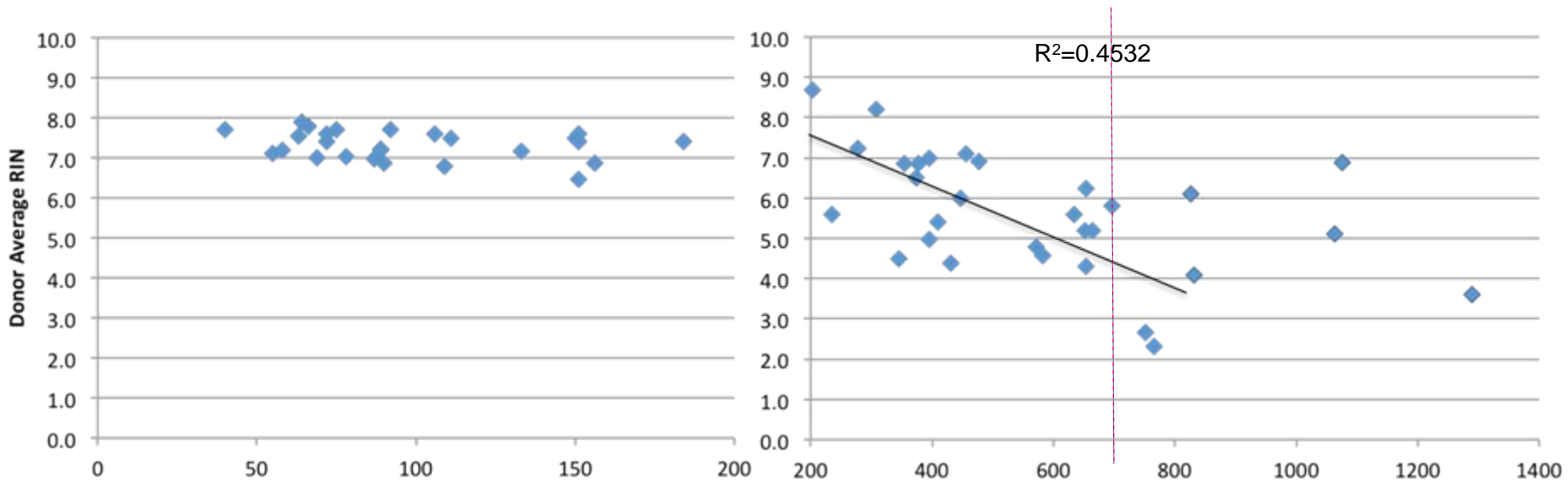
RNA Quality and PMI



RIN decreases with increasing PMI for PAXgene Tissues, but not for Fresh Frozen Brains

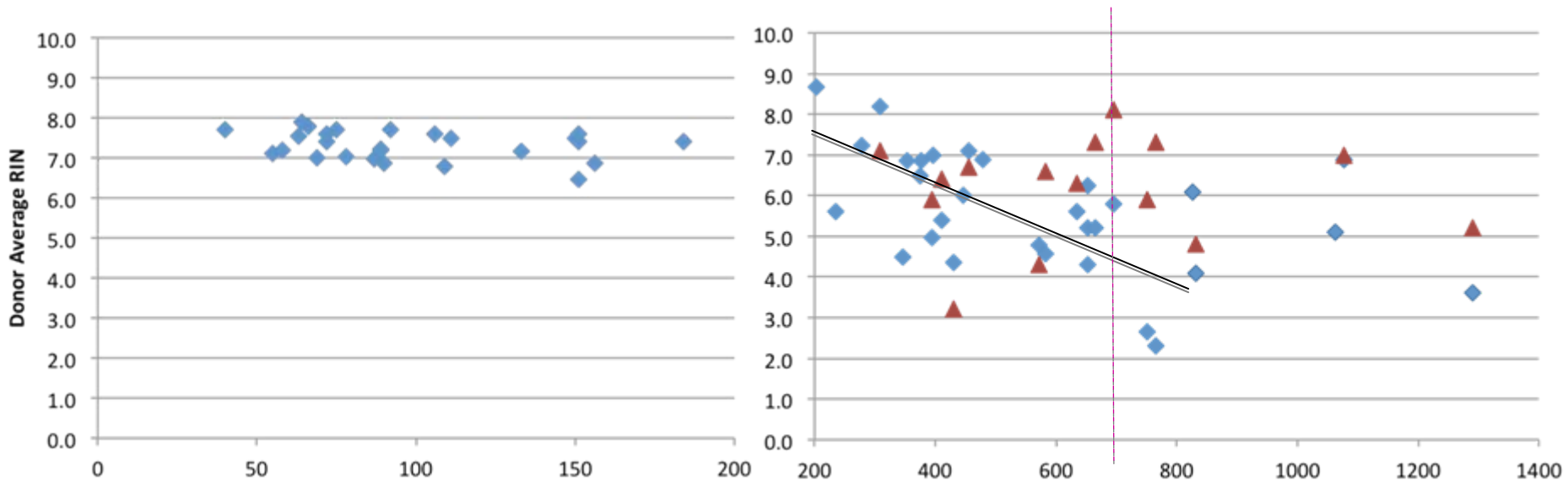


RNA Quality declines with PMI



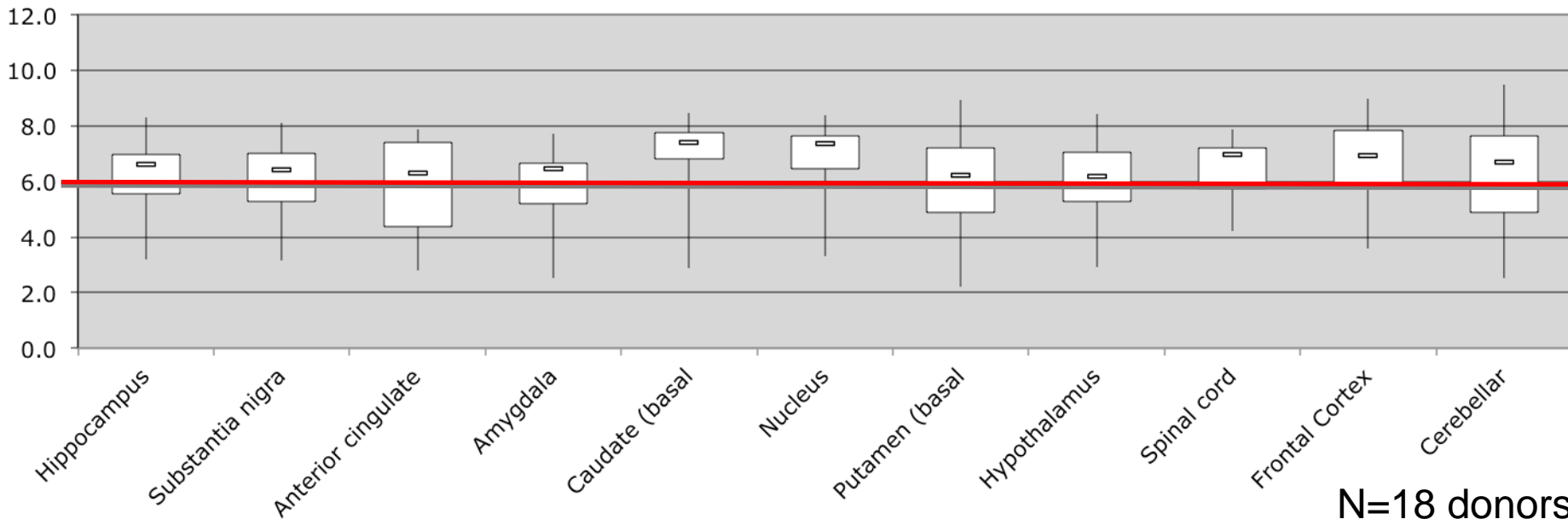
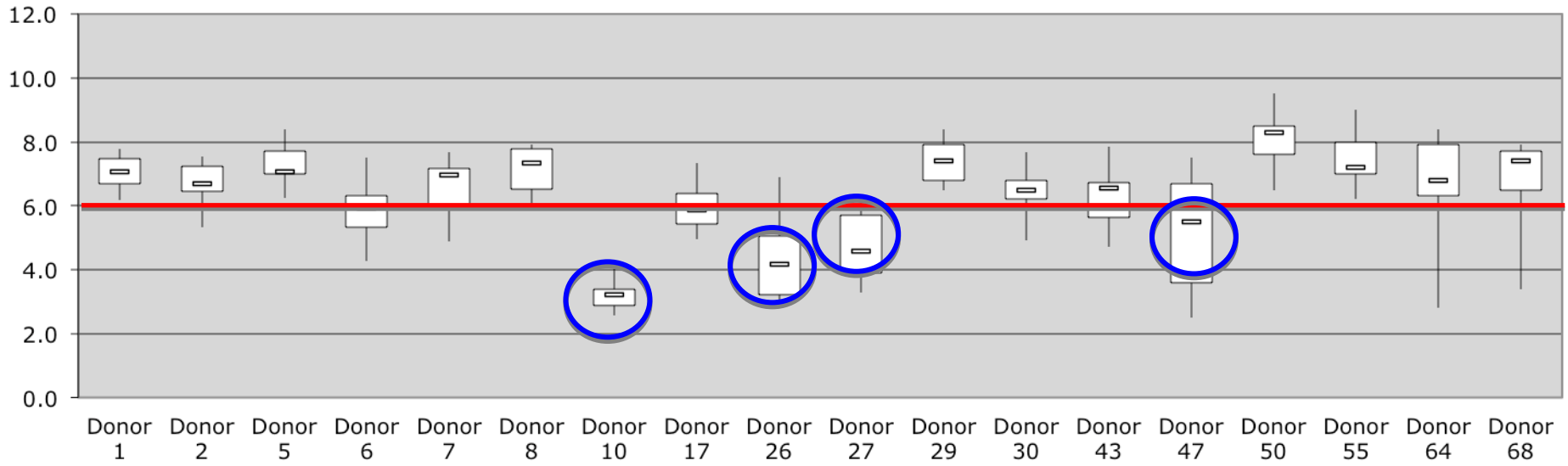
RINs are largely stable out to about 3 hours, then decline from 3- ~12 hours. More variable thereafter.

RNA Quality declines with PMI



Brain RINs less correlated with PMI, and not always the same as mean PAXgene RIN from same donor

Fresh Frozen Brains

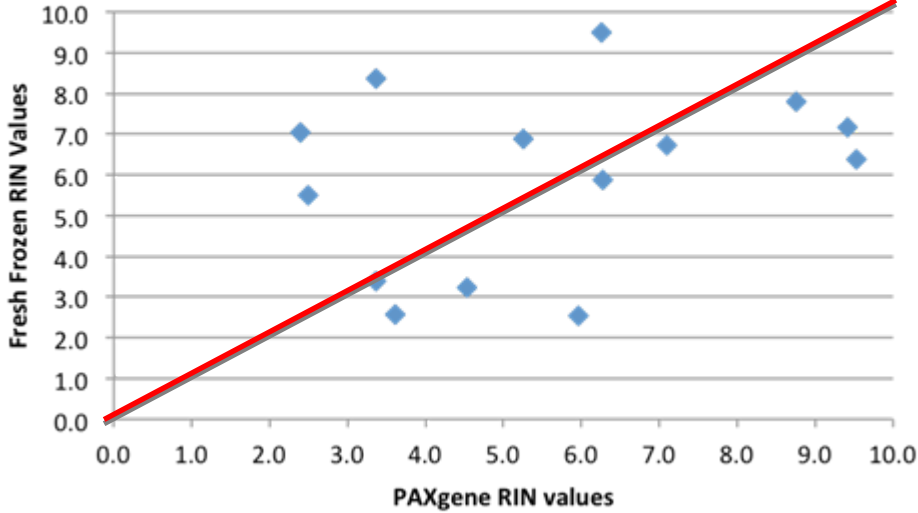


N=18 donors

PAXgene vs FF RNA Quality - Brains

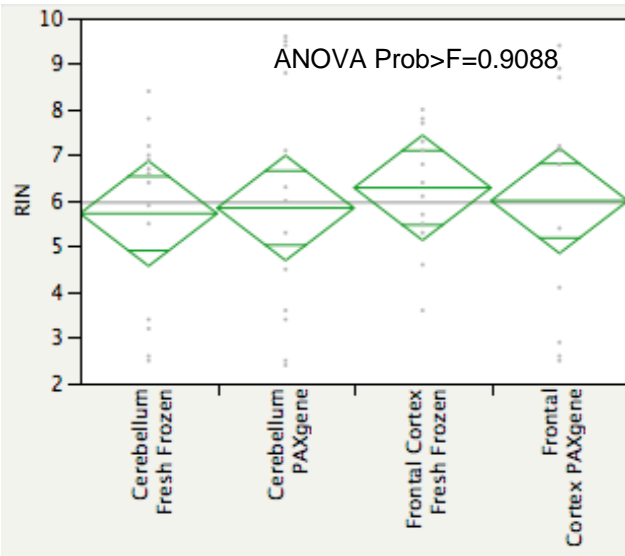
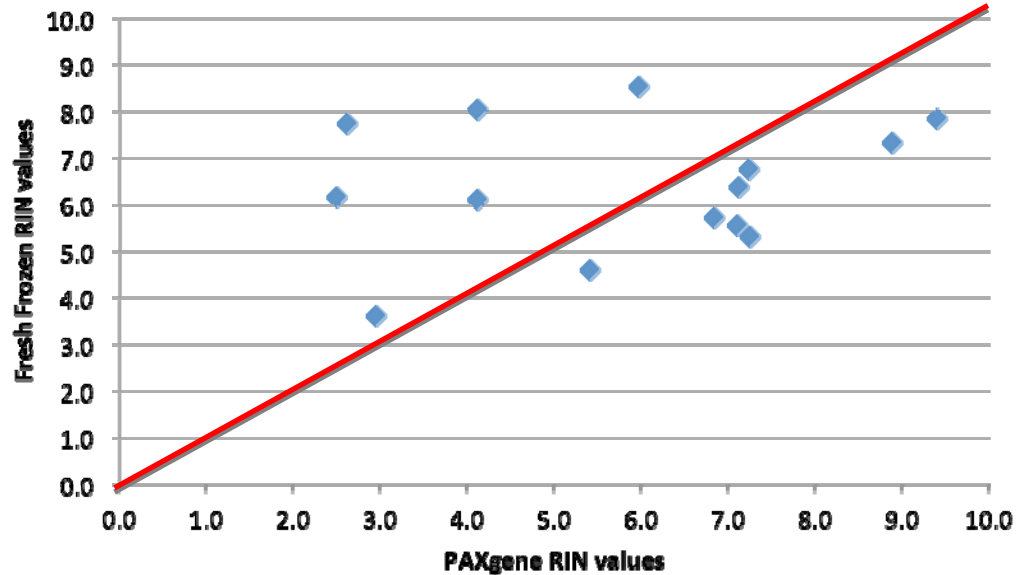


Cerebellum

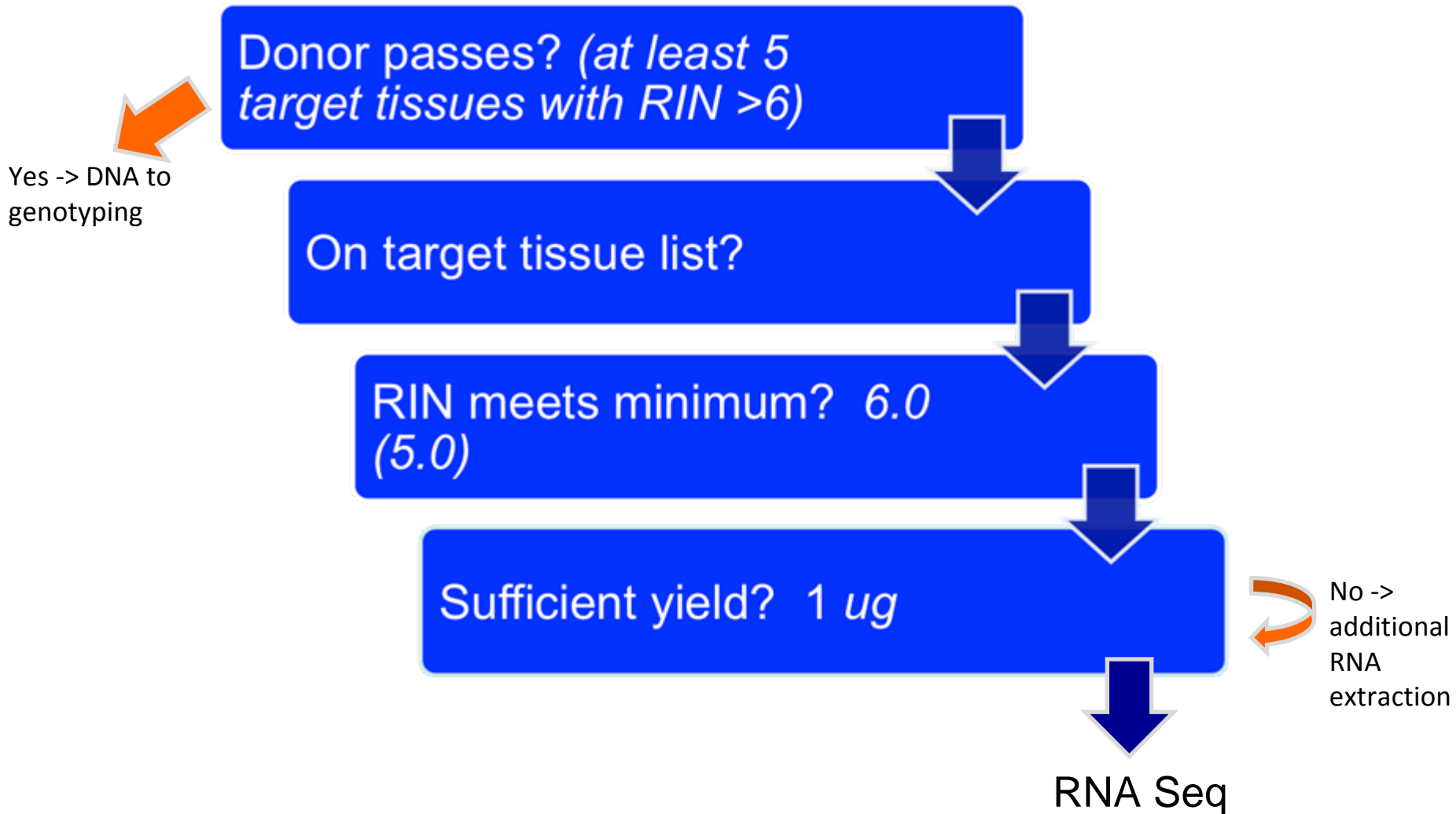


N=15 donors

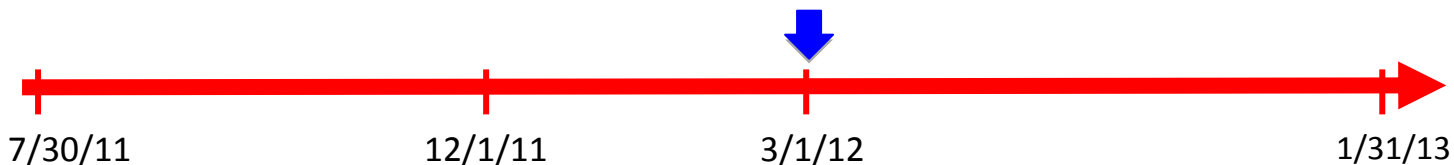
Cortex



Sample Qualification



Molecular Analysis Current Status



Genotyping Arrays



Expression Arrays



RNA-Seq

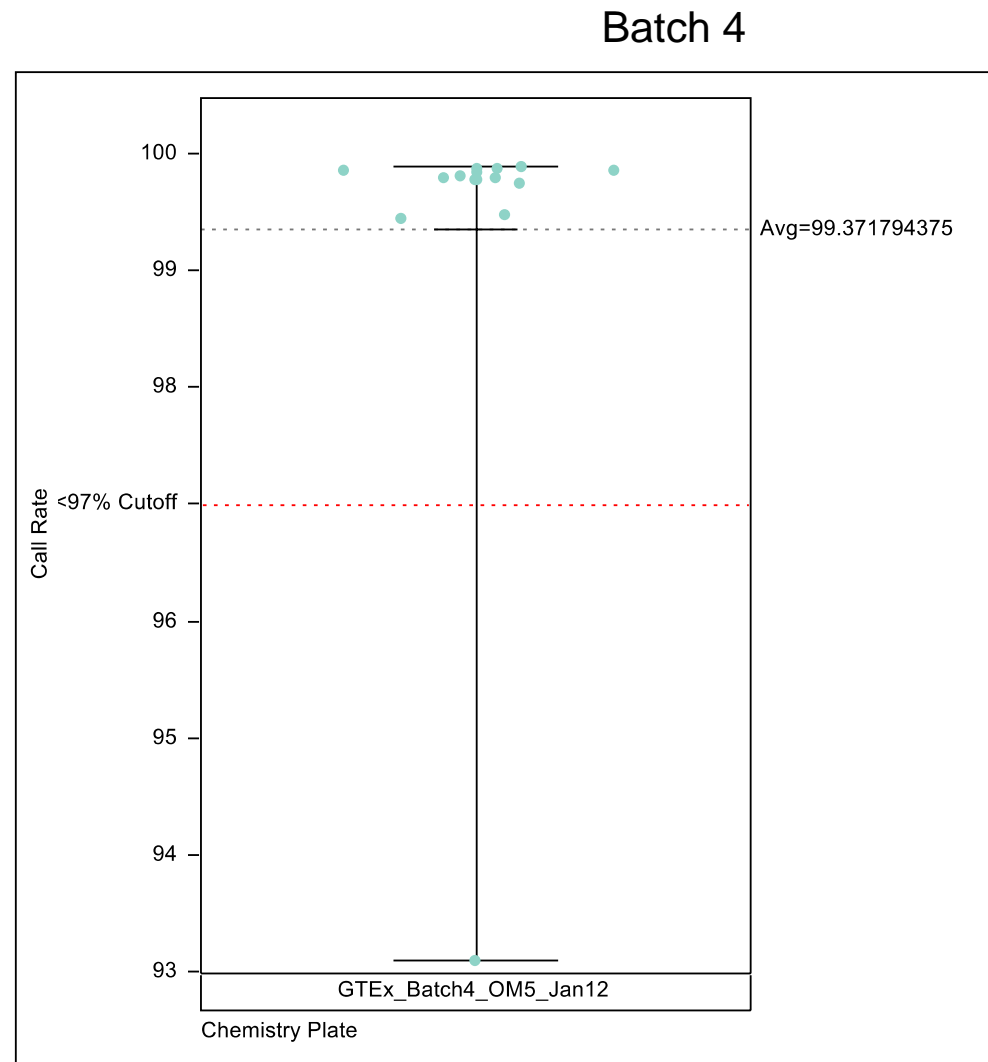


Completed by 3/1/12 Data freeze	Pilot Goal
64	300
376 (34)	1000
290 (38)	1500

Illumina 5M Genotyping

- Allele Frequency Distribution of ~5 million SNPs in GTEx 5M Batch 1
 - 100% pass rate
 - 99.35% average call rate
 - ~3.1 million polymorphic sites

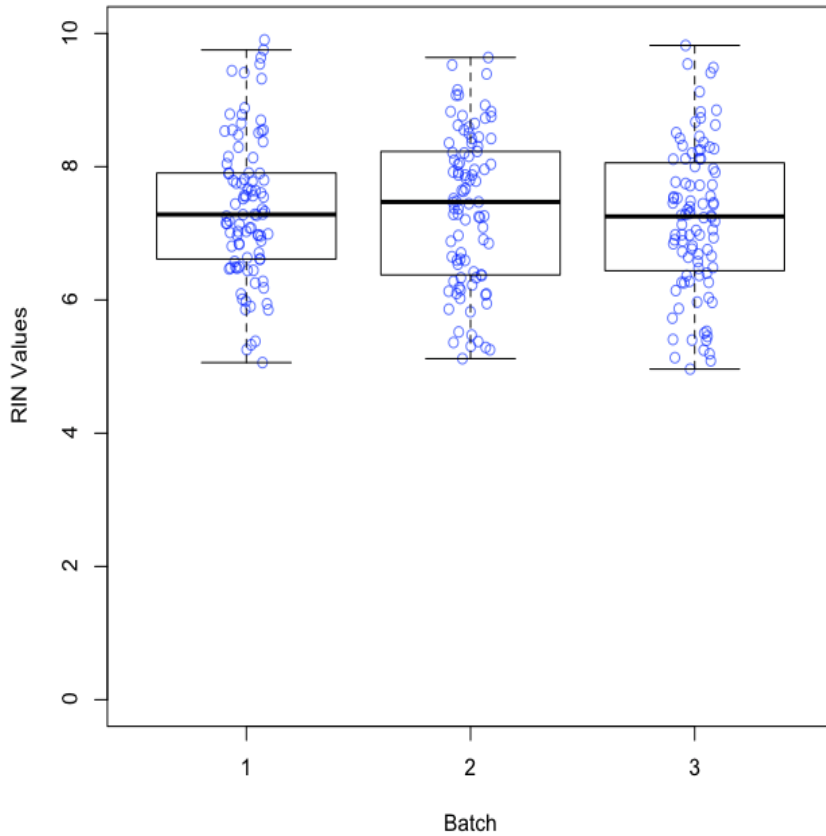
MAF	SNP Count
0	1156894
0.03-0.1	1615589
0.1-0.2	530866
0.2-0.3	362103
0.3-0.4	337539
0.4-0.5	298341



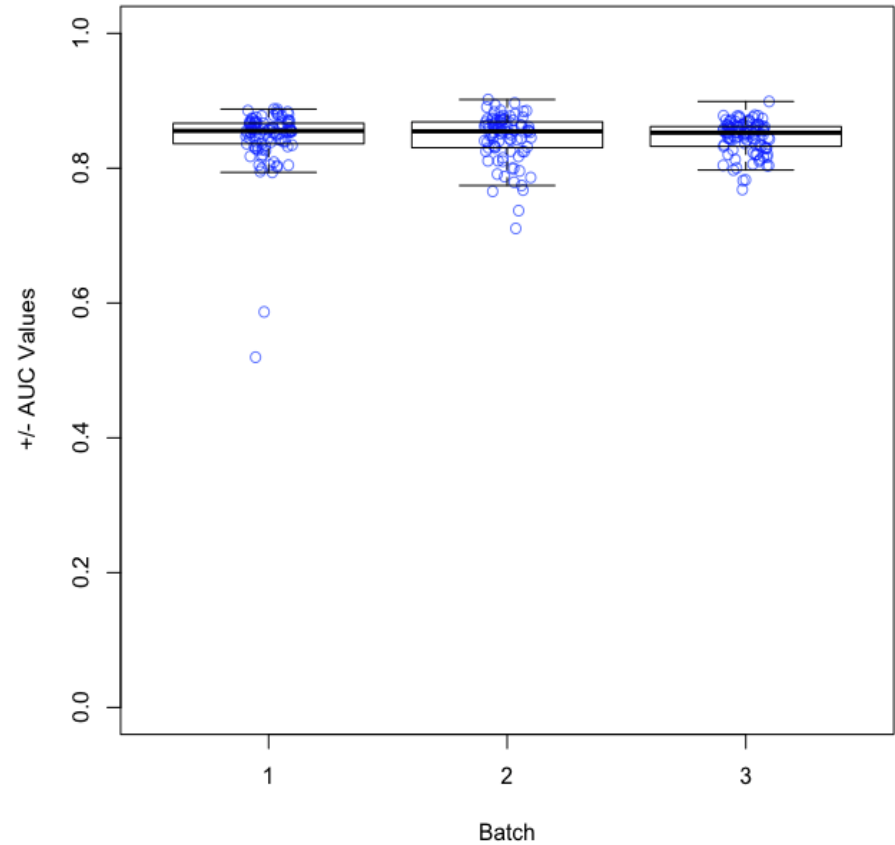
Expression Array QC by Batch



GTEx RINs by Batch

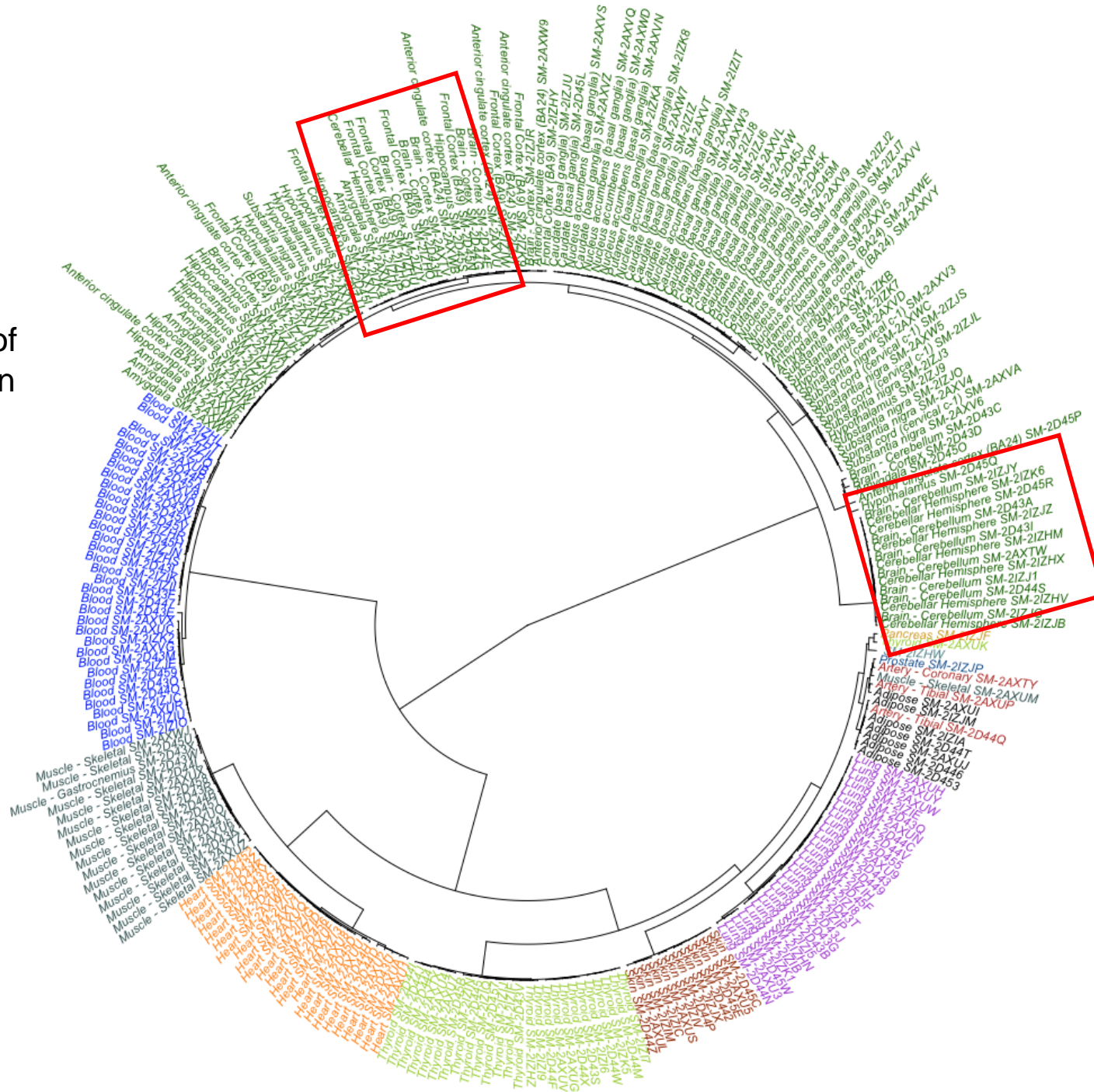


GTEx +/- AUC by Batch



Hierarchical Clustering

Distance measure:
Spearman Correlation of
all Affymetrix expression
levels



RNA Sequencing

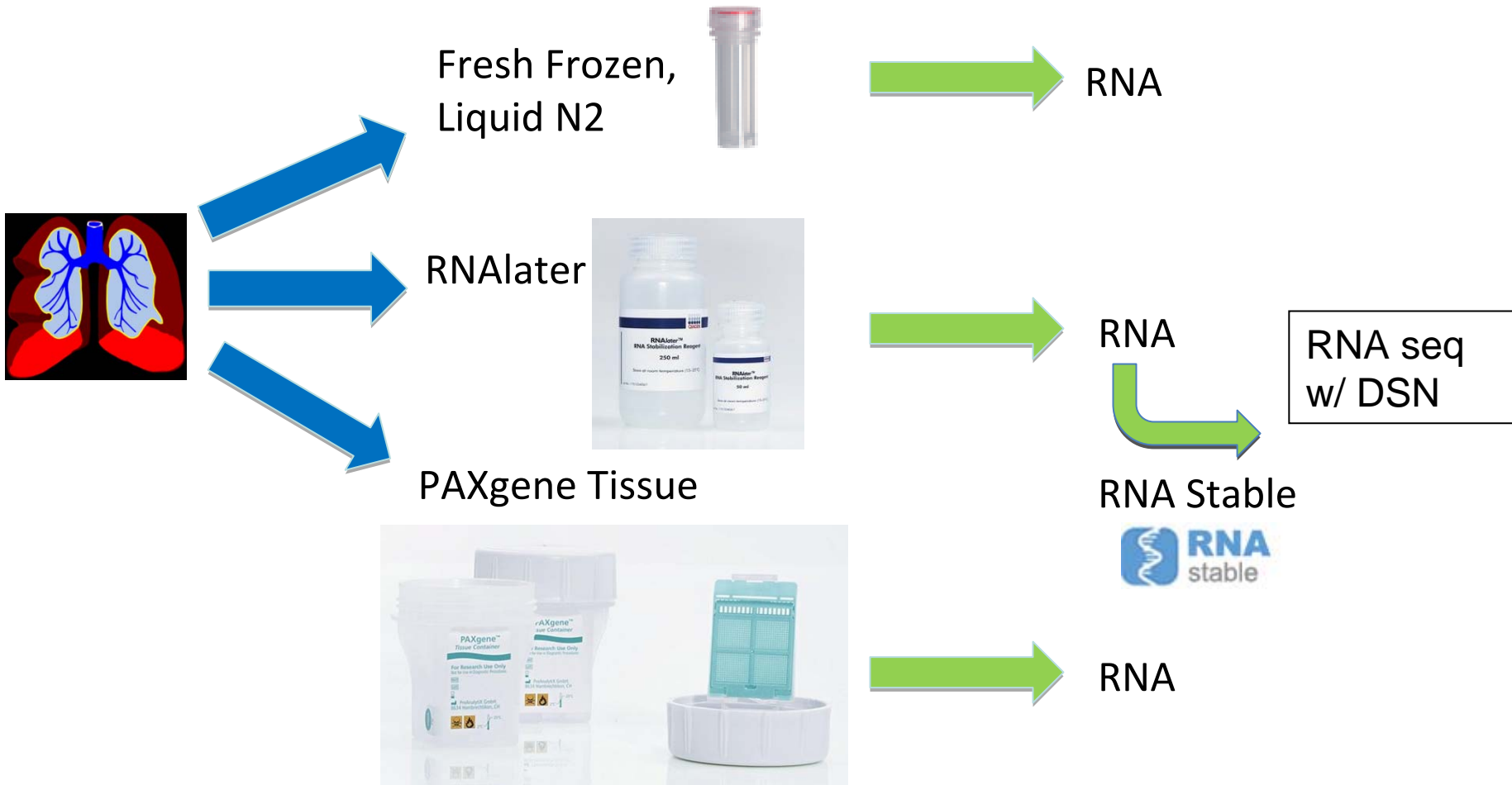


- Tech dev efforts have focused on identifying a cDNA library construction method that provides high quality transcript profiles and:
 - Uses < 1ug total RNA input
 - Works with degraded RNA
 - Can be readily scaled to high-throughput
 - Works with preserved tissues (e.g. Paxgene Tissue)

Poly-A method (Tru-Seq, dUTP) vs non-poly-A DSN

GTEx Pilot RNA Experiments

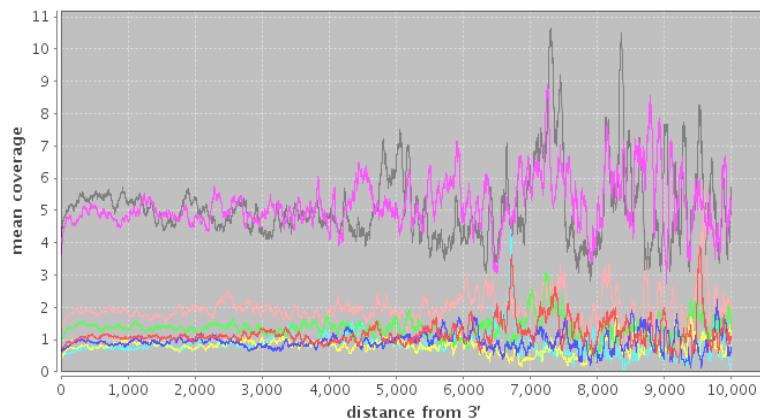
GOAL → Select tissues for which the 3 methods produced similar, and at least mid-range, RIN scores for all samples, and for which there was sufficient RNA yield for both Affymetrix GeneArrays and RNASeq experiments.



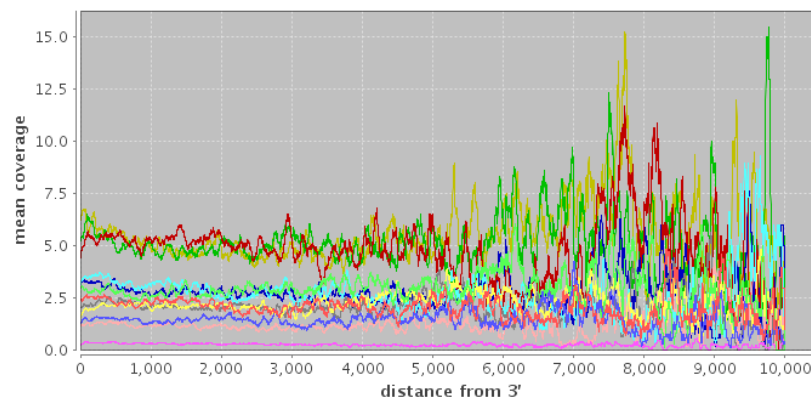
GTEEx Pilot Tissue experiment results



*Based on 1000 middle expressed transcripts



GTEX-1663F1C8-0503 GTEX-1663F1C8-0520_Sample1 GTEX-1663F1C8-0520_Sample2
 GTEX-1663F1C8-2619 GTEX-1663F1C8-0527 GTEX-1663F1C8-2604 GTEX-1663F1C8-2627
 GTEX_K_33rdDSN_an8



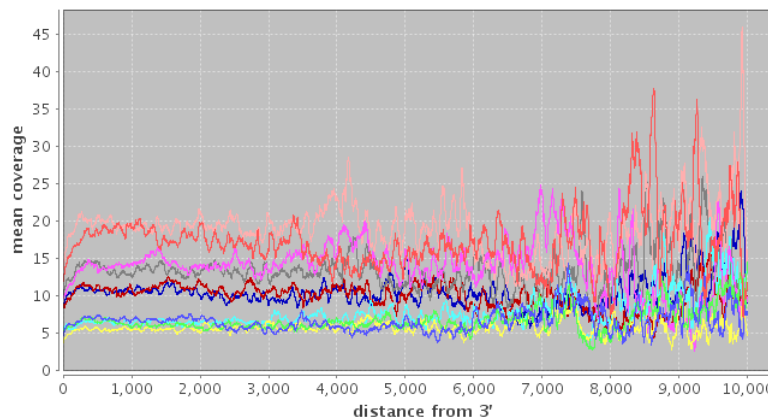
GTEX-1809447E-0726 GTEX-46415CE4-0501 GTEX-46415CE4-0517 GTEX-46415CE4-0525
 GTEX-46415CE4-1201 GTEX-46415CE4-1217 GTEX-46415CE4-1225_sample1
 GTEX-46415CE4-1225_sample2 GTEX-46415CE4-1801 GTEX-46415CE4-1817
 GTEX-46415CE4-1825 K562_Total_RNA_Ambion

Donor 2

Mean dup. rate: 73%
 Mean CV: 1.23
 Mean Gap %: 41%

Donor 3

Mean dup. rate: 54%
 Mean CV: 0.93
 Mean Gap %: 11%



GTEX_35thDSN_K562 GTEX-1809447E-0701 GTEX-1809447E-0717 GTEX-1809447E-0801
 GTEX-1809447E-0817 GTEX-1809447E-0826 GTEX-1809447E-1501_sample1
 GTEX-1809447E-1501_sample2 GTEX-1809447E-1517 GTEX-1809447E-1526

Donor 4

Mean dup. rate: 64%
 Mean CV: 1.00
 Mean Gap %: 19%

**No evident
 difference
 between
 preservation
 methods.**

RNA Seq First Production Batch



Donor	Tissue	RIN
GTEX-N7MS	Blood	8.2
GTEX-N7MS	Brain	8.7
GTEX-N7MS	Brain	9.6
GTEX-N7MT	Lung	9.1
GTEX-N7MT	Brain	9.5
GTEX-NPJ8	Blood	5.5
GTEX-NPJ8	Lung	8.8
GTEX-NPJ8	Brain	8.8
GTEX-NPJ8	Brain	9.4
GTEX-O5YV	Lung	6.1
GTEX-O5YV	Muscle	7.1



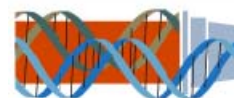
4 ug polyA+ LC
(2 methods)



1 ug DSN LC



1/6th lane, 101 PE

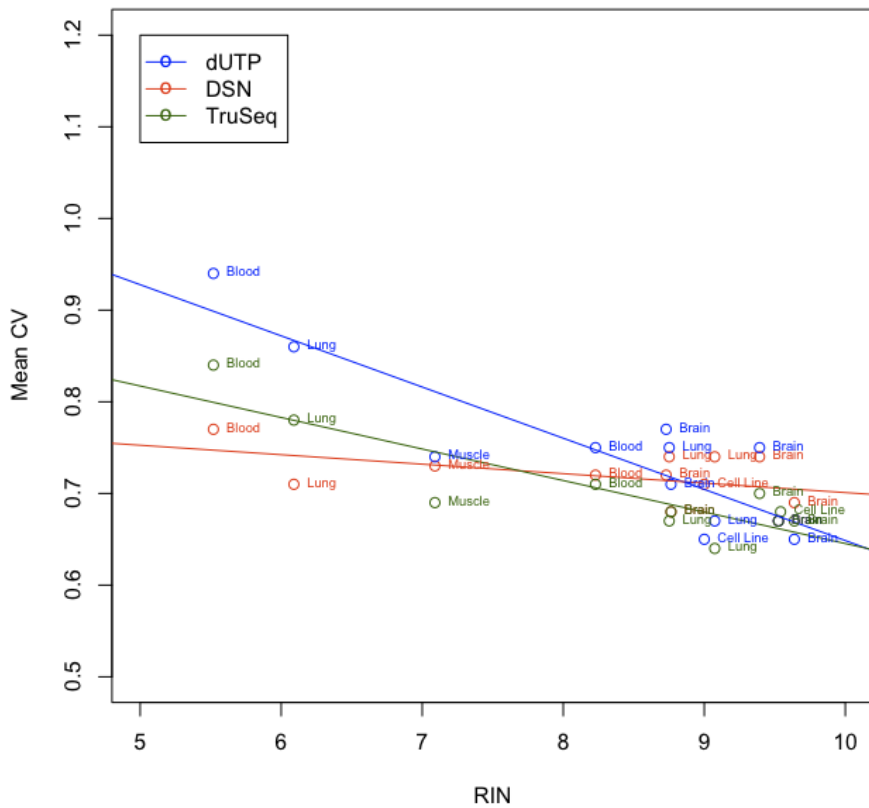


FIREHOSE
Genotype-Tissue Expression

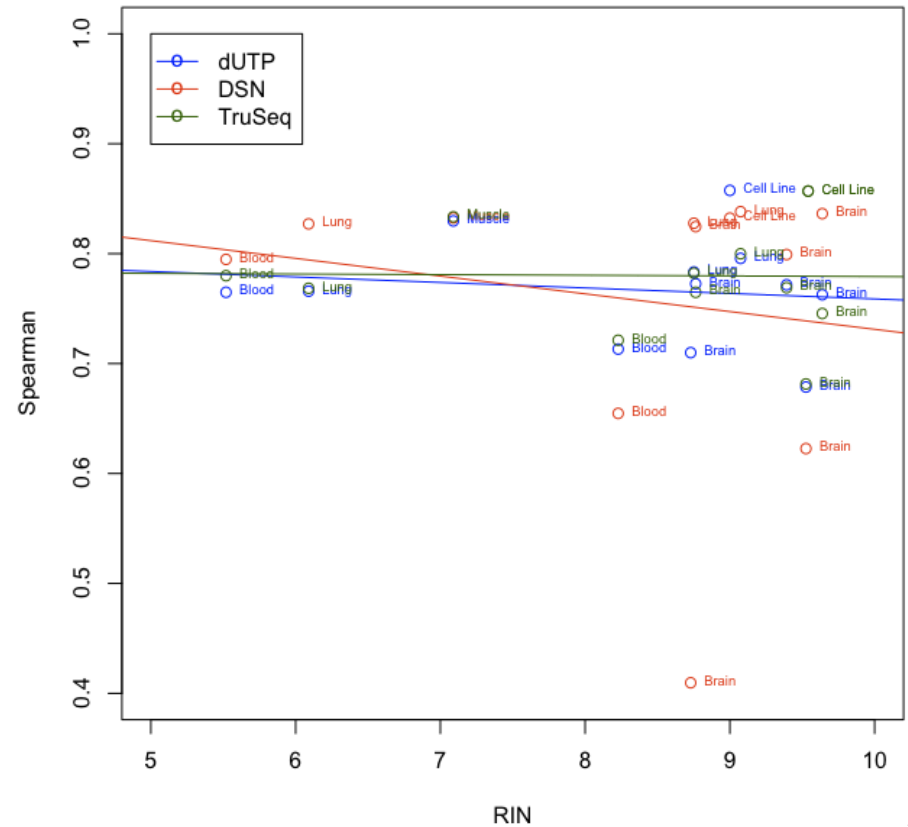
Coverage metrics



Effect of RIN on Mean CV



Effect of RIN on Spearman



For middle expressed 1000 transcripts, downsampled to 20m reads.

LC Protocol Summary



- For RINs ≥ 7 , polyA+ and DSN perform similarly
 - polyA does seem to be adversely affected by RINs < 6
- BUT, DSN has significant disadvantages compared to polyA+ in terms of protocol robustness, GC bias, and cost
- polyA+ protocol currently requires 1-5ug RNA, but ongoing testing aims to reduce that to 0.1-1ug

CURRENT PRODUCTION – polyA+ protocol, with 1ug RNA input and RIN ≥ 6

RNA Seq Analysis Pipeline

firehose:8080/gtex/secure/annotatedentity/report/Sample/12728302/?branchId=3004

V2.1.5

FIREHOSE
Genotype-Tissue Expression

Wednesday, November 16, 2011 12:38:02 AM EST
Logged in as : ddeleuca [logout](#)
Search Entity

Workspaces ▾ Samples ▾ Individuals ▾ Pipelines ▾ Admin ▾ Help

Workspace: Tech-Dev

REPORTS FOR SAMPLE [GTEx-1663F1C8-0520_SAMPLE1](#)

Show/Hide All

RNA-SEQ QC METRICS

RNA-SEQ METRICS

READ COUNT METRICS

The following summary statistics are calculated by counting the number of reads that have the given characteristics.

TOTAL READS

Sample	Note	Total	Unique	Duplicates	Duplication Rate	Estimated Library Size
GTEx-1663F1C8-0520_Sample1	Thyroid - Fresh Frozen	32,187,140	8,507,110	23,680,030	0.736	1,000,574

Total reads are filtered for vendor fail flags. **Unique** are reads without the duplicate flag. **Duplicates** are reads with duplicate flag. **Duplication Rate** is the number of duplicate reads divided by total reads. **Estimated Library Size** is the number of expected fragments based upon the total number of reads and duplication rate assuming a Poisson distribution.

MAPPED READS

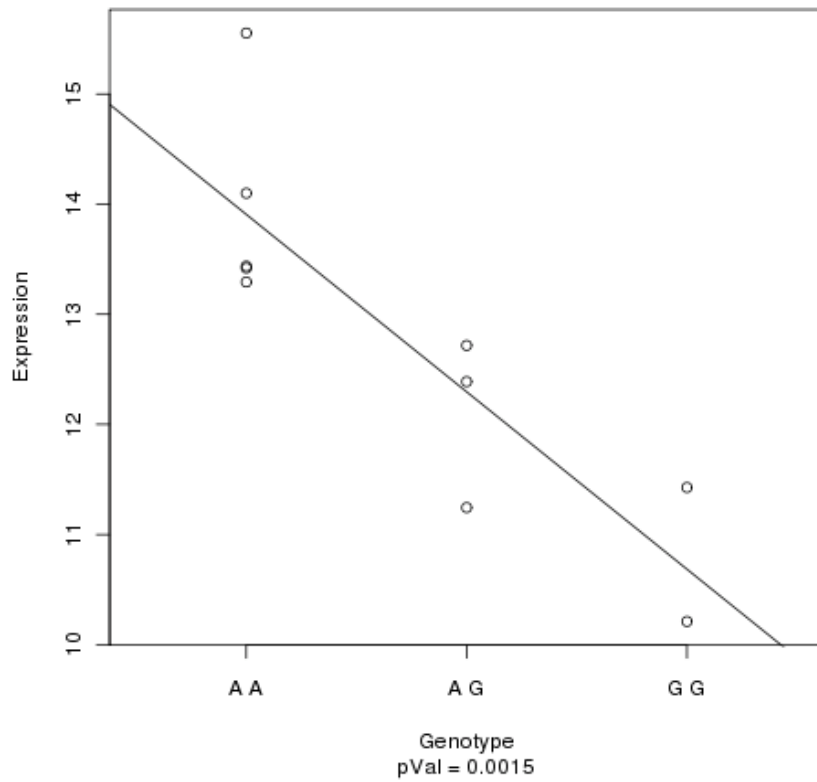
Sample	Note	Mapped	Mapping Rate	Mapped Unique	Mapped Unique Rate	rRNA	rRNA rate	Unique non-rRNA	Unique non-rRNA Rate
GTEx-1663F1C8-0520_Sample1	Thyroid - Fresh Frozen	26,084,555	0.810	2,404,525	0.075	9,068,987	0.282	-6,664,462	-0.207

Mapped reads are those that were aligned. **Mapping Rate** is per total reads. **Mapped Unique** are both aligned as well as non-duplicate reads. **Mapped Unique Rate** is per total reads. **rRNA** reads are non-duplicate and duplicate reads aligning to rRNA regions as defined in the transcript model definition. **rRNA Rate** and **Unique non-rRNA Rate** are per total reads.

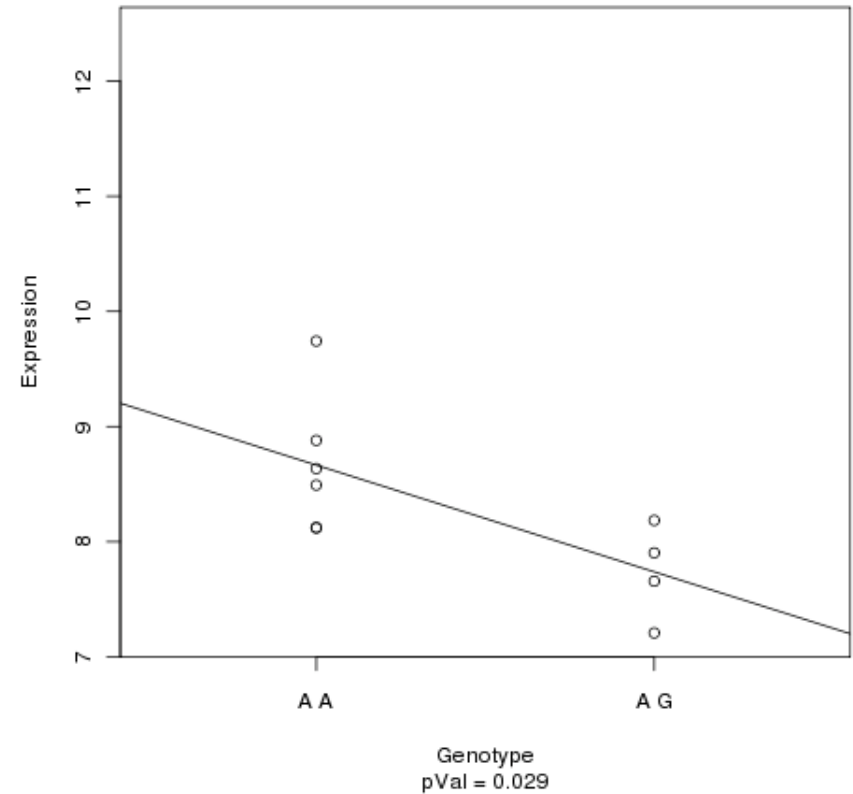
Known Liver eQTLs found in Lung



rs599839 SORT1 Lung



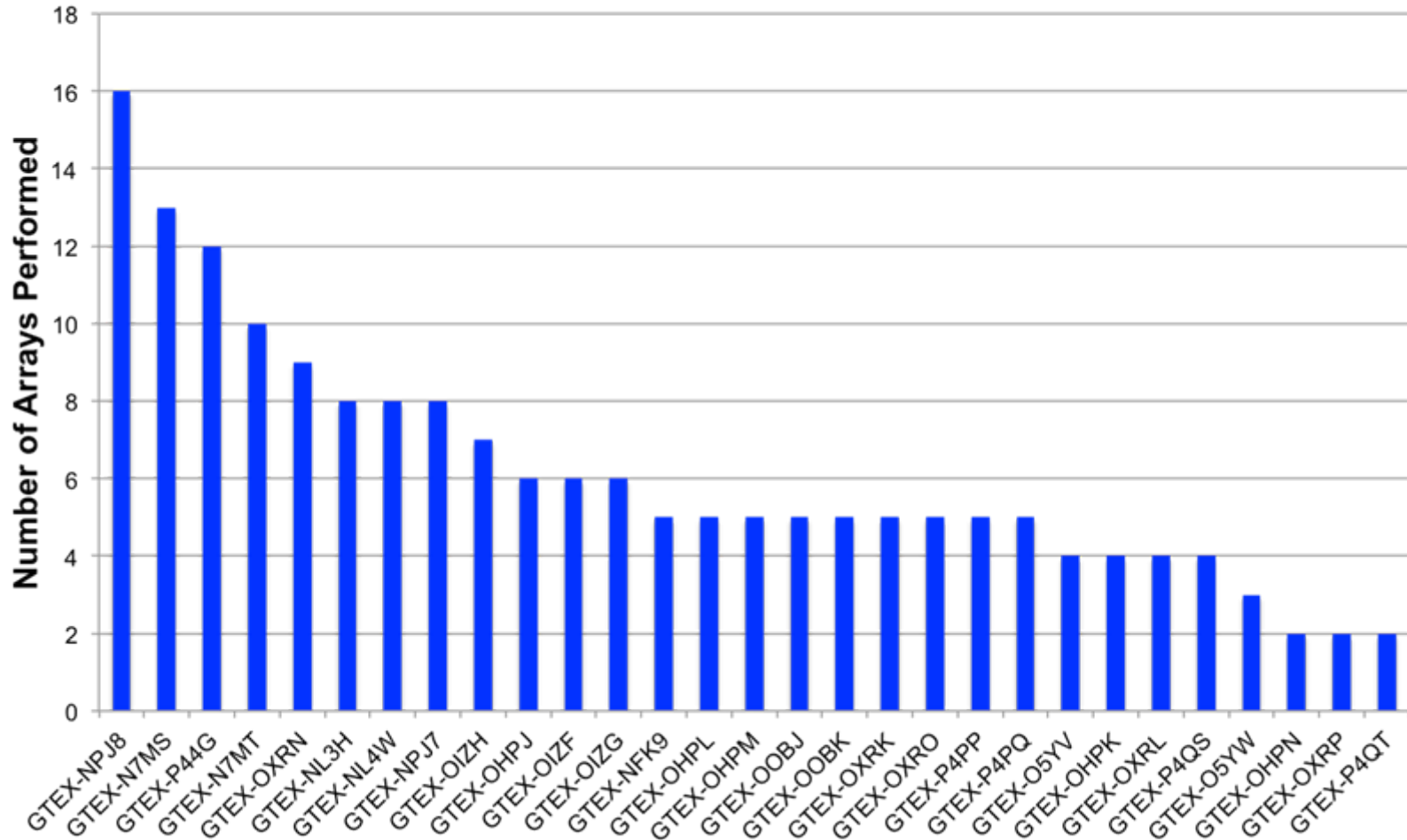
rs7013657 XKR9 Lung



Current Power Limitations



Expression Arrays by Donor



Broad LDACC Acknowledgements



BSP

Julie Ann
Clint Chalk
Cara Fischer
Kim Kanki
Tsheko Makuwa
Paula Morais
BSP Staff
Kristin Ardlie

INFORMATICS

Michael Dinsmore
George Grant
Alex Thomson
Taylor Young
Tim Fennell
GAP, BSP, GSP teams

GAP

Yves Boie
Andrew Crenshaw
Supriya Gupta
Melissa Parkin
Prapti Pokharel
GAP Staff
Wendy Winckler (co-PI)

PROJECT MANAGEMENT

Molly Donovan
Ellen Gelfand
Kristin Thompson

ANALYSIS

David Deluca
Andrey Sivachenko
Gad Getz

GSP

Sheila Fisher
Ryan Johnson
Paula Maness
Jane Wilkinson
GSP staff
Jen Baldwin

GSAP

Xian Adiconis
Jim Bochicchio
Joshua Levin
Carsten Russ
Chad Nusbaum

Harvard Skin Disease Research Center

Jim Rheinwald
Patty Barron