

RNA editing

New sources of genomic variation, software and application

Scott Funkhouser

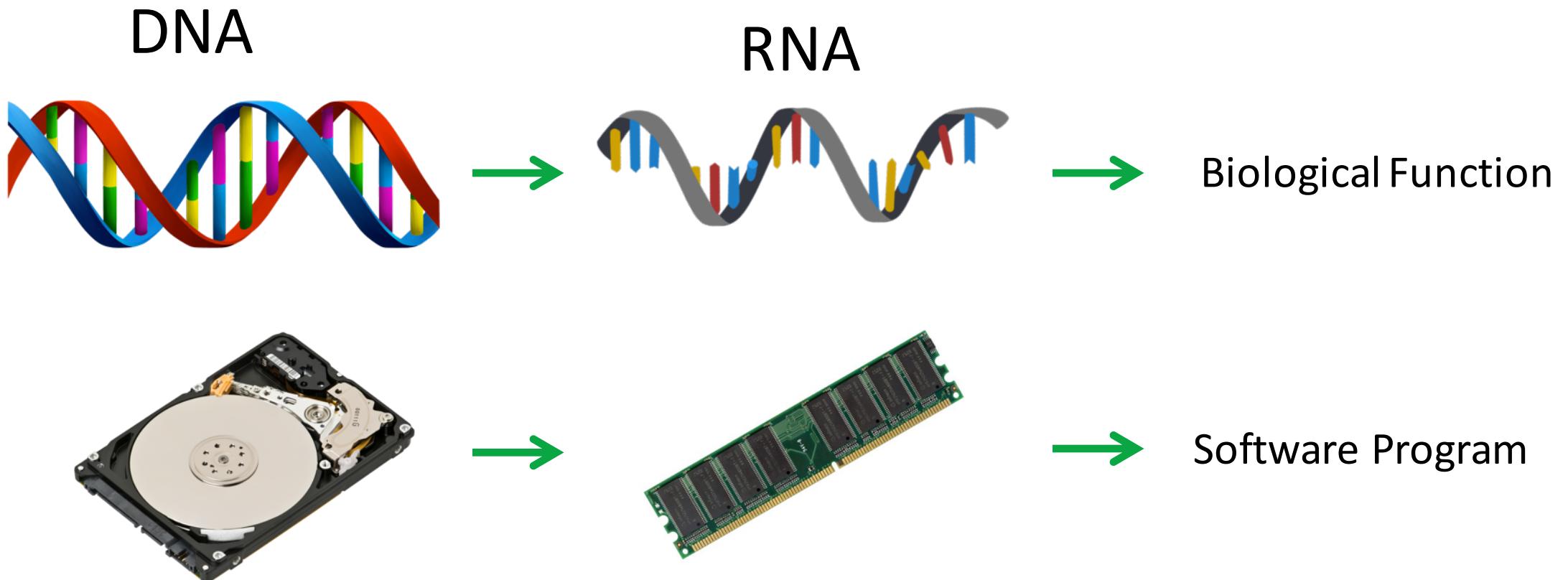
Ernst Lab

Genetics and Quantitative Biology PhD Candidate

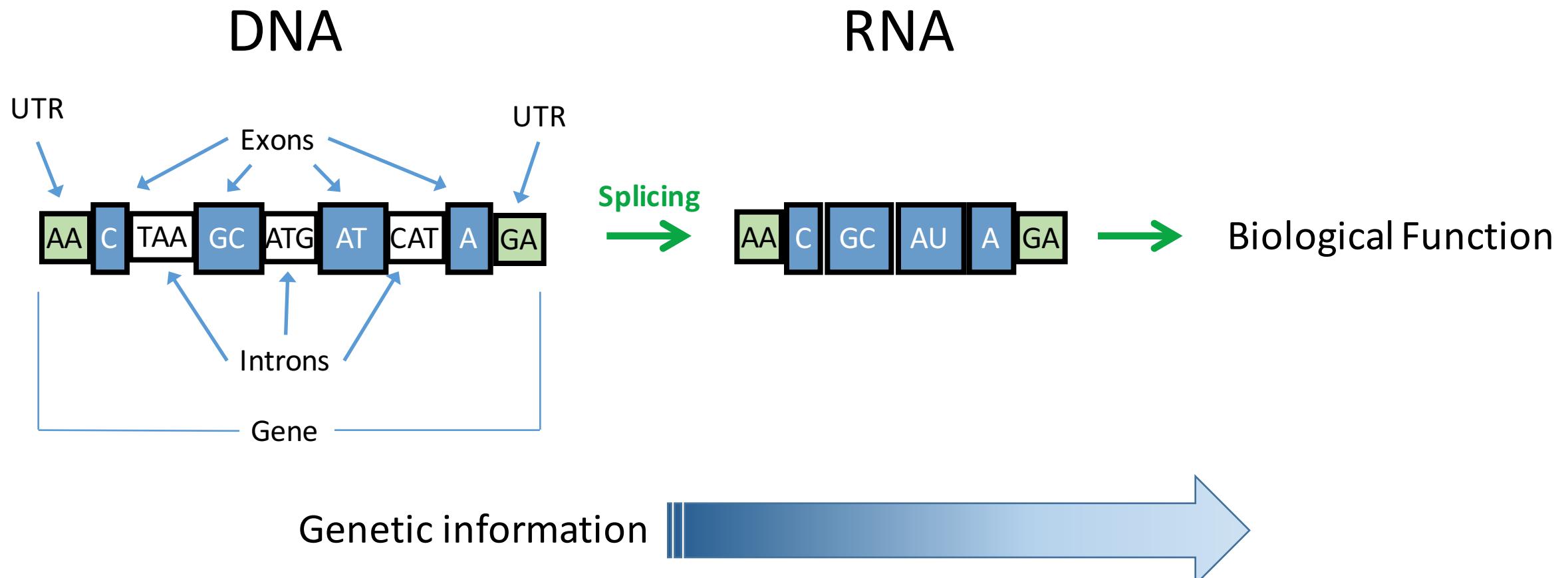
Outline

- What is RNA editing?
- *editTools*: an R package for RNA editing analysis
- Application of *editTools* to the swine model

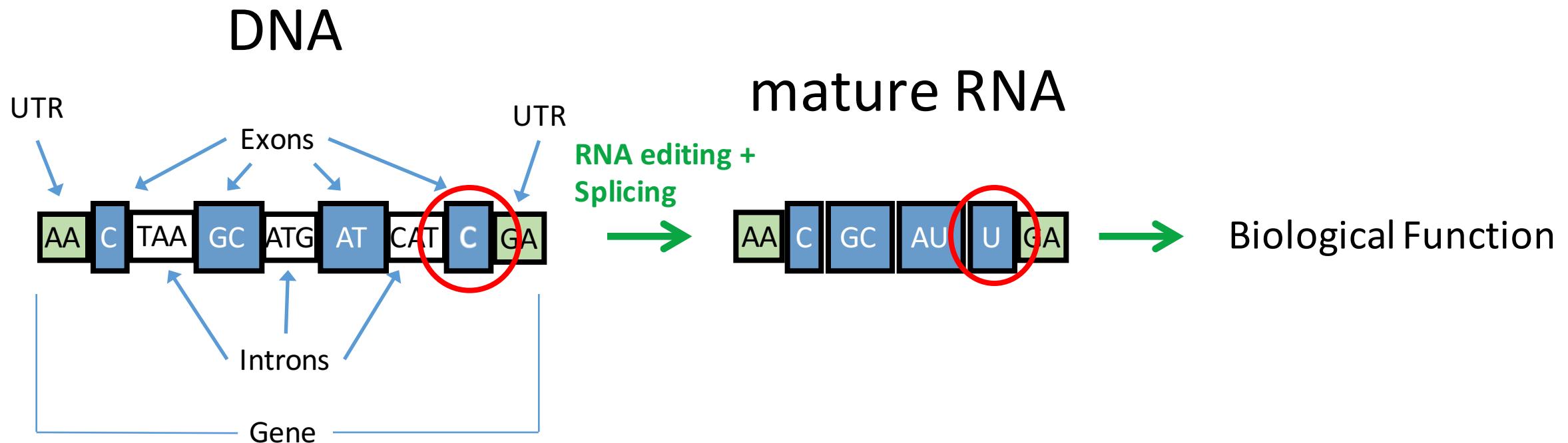
The central dogma



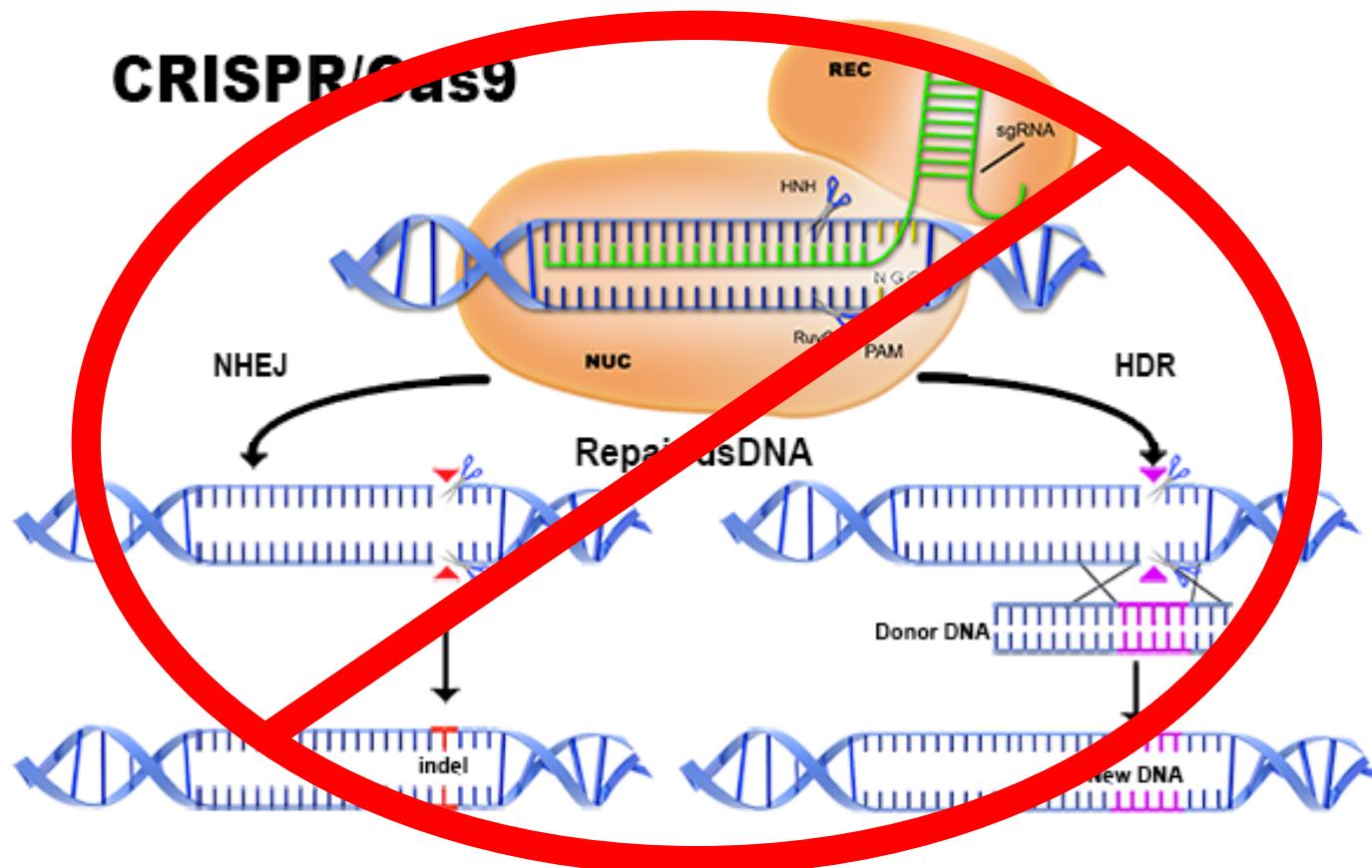
The central dogma



RNA editing



Not to be confused with “genome editing”



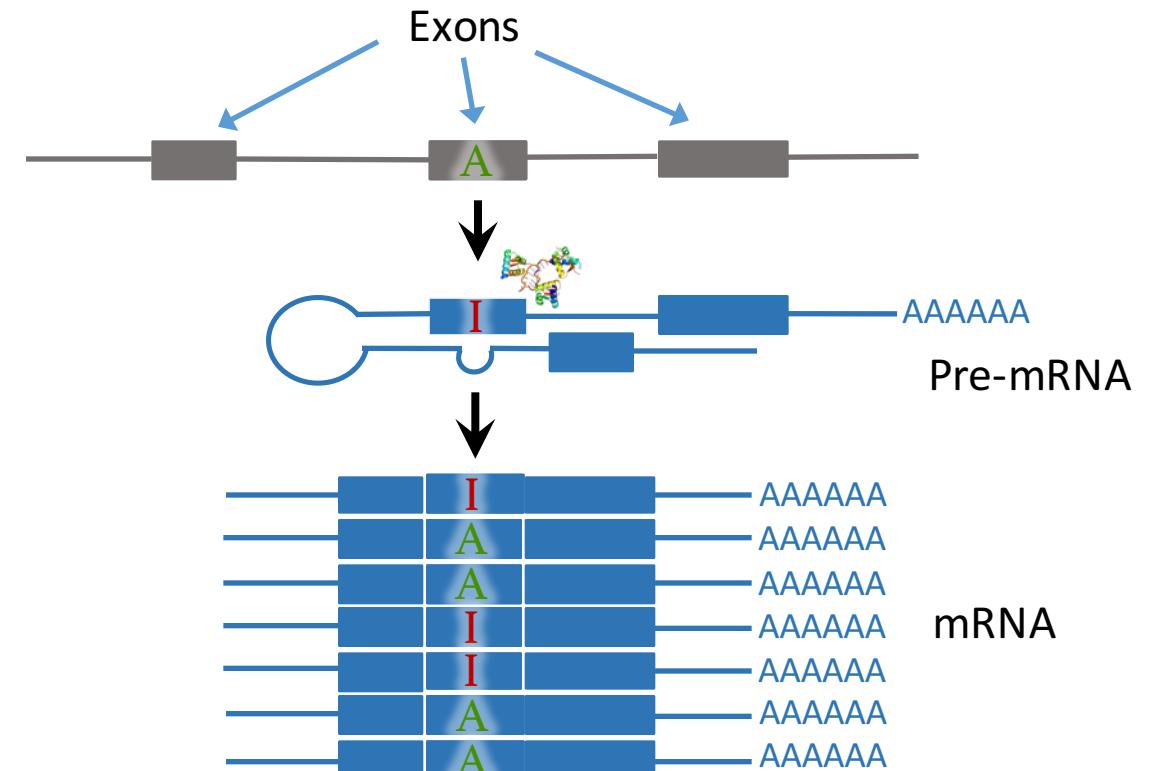
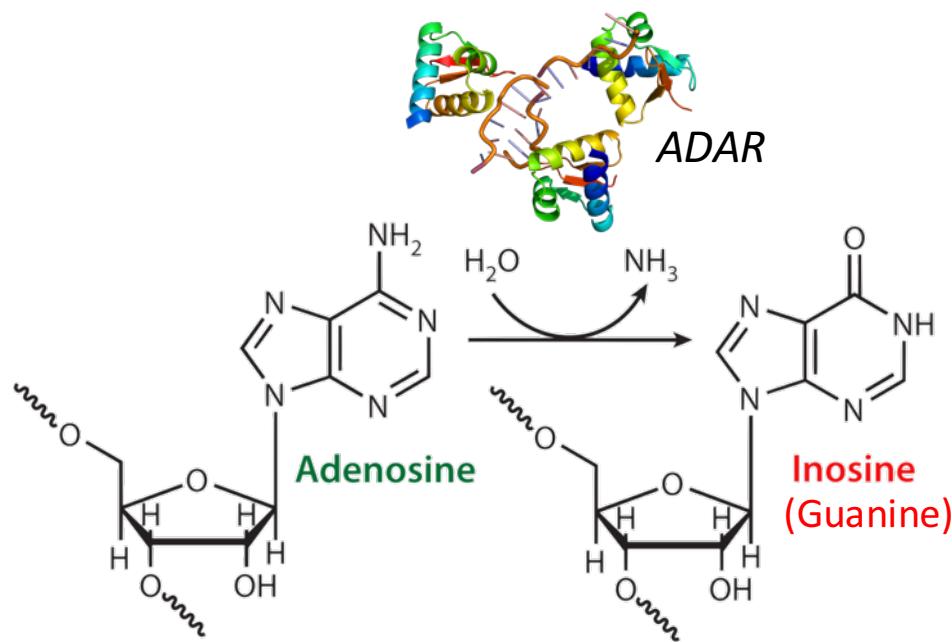
Genome editing by CRISPR cas9

- Editing performed on the genome
- Artificial technology

RNA editing

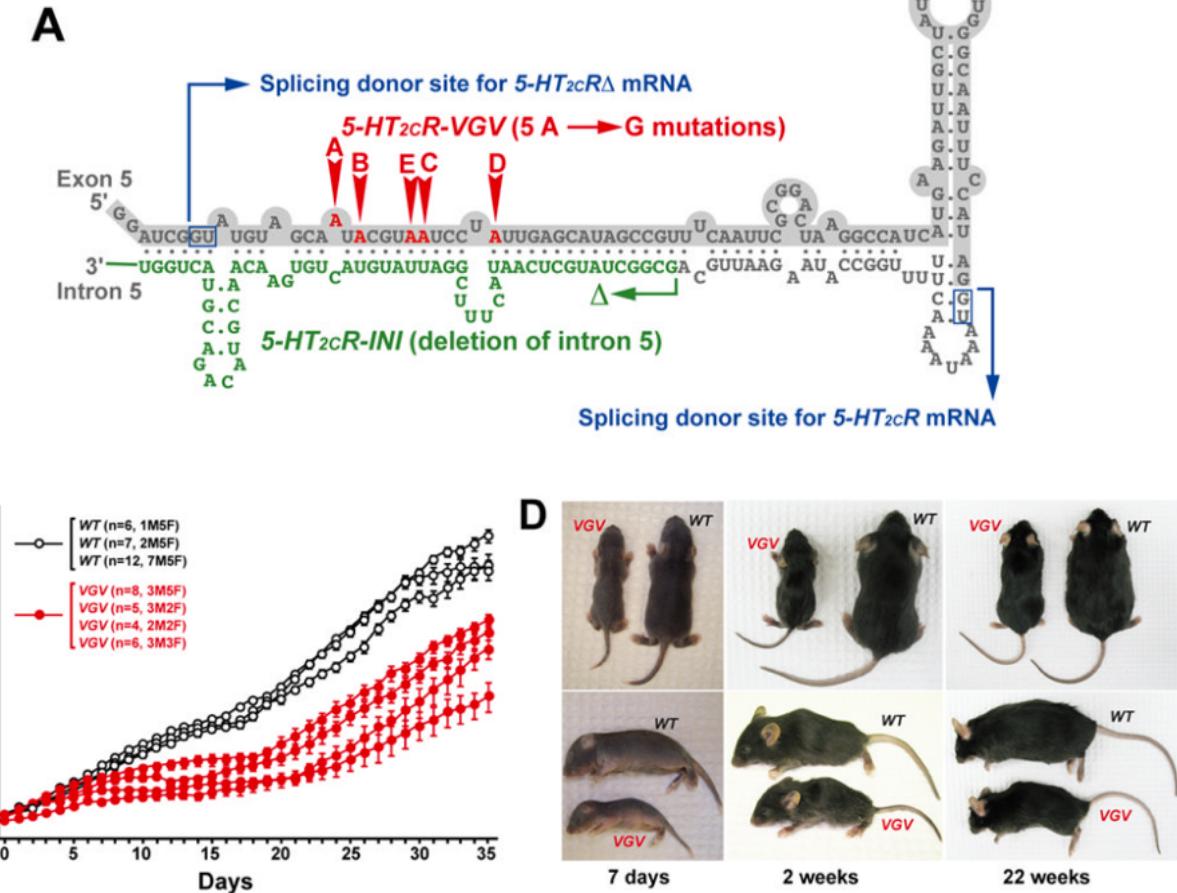
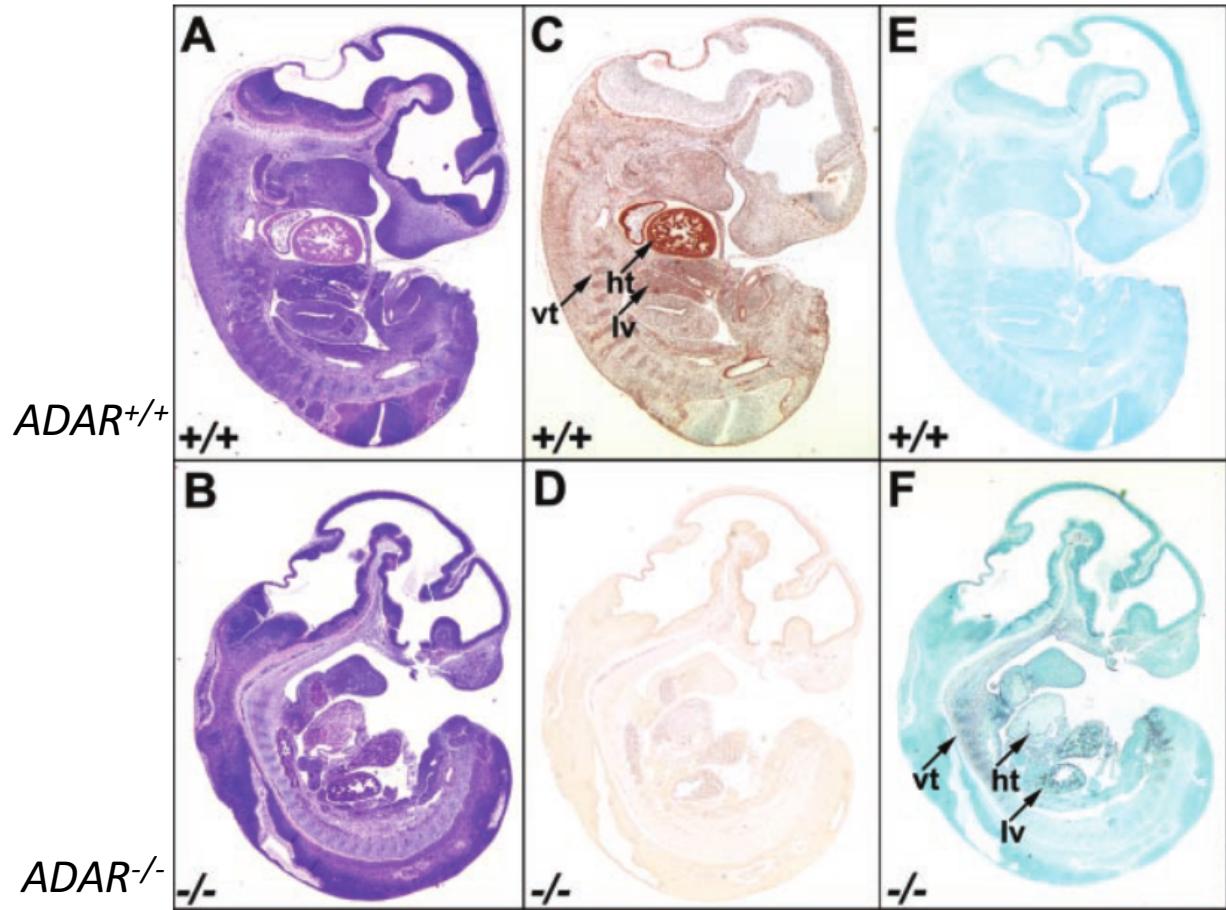
- Editing performed on the transcriptome
- Natural. Done intentionally by the cell.

RNA editing in mammalian models



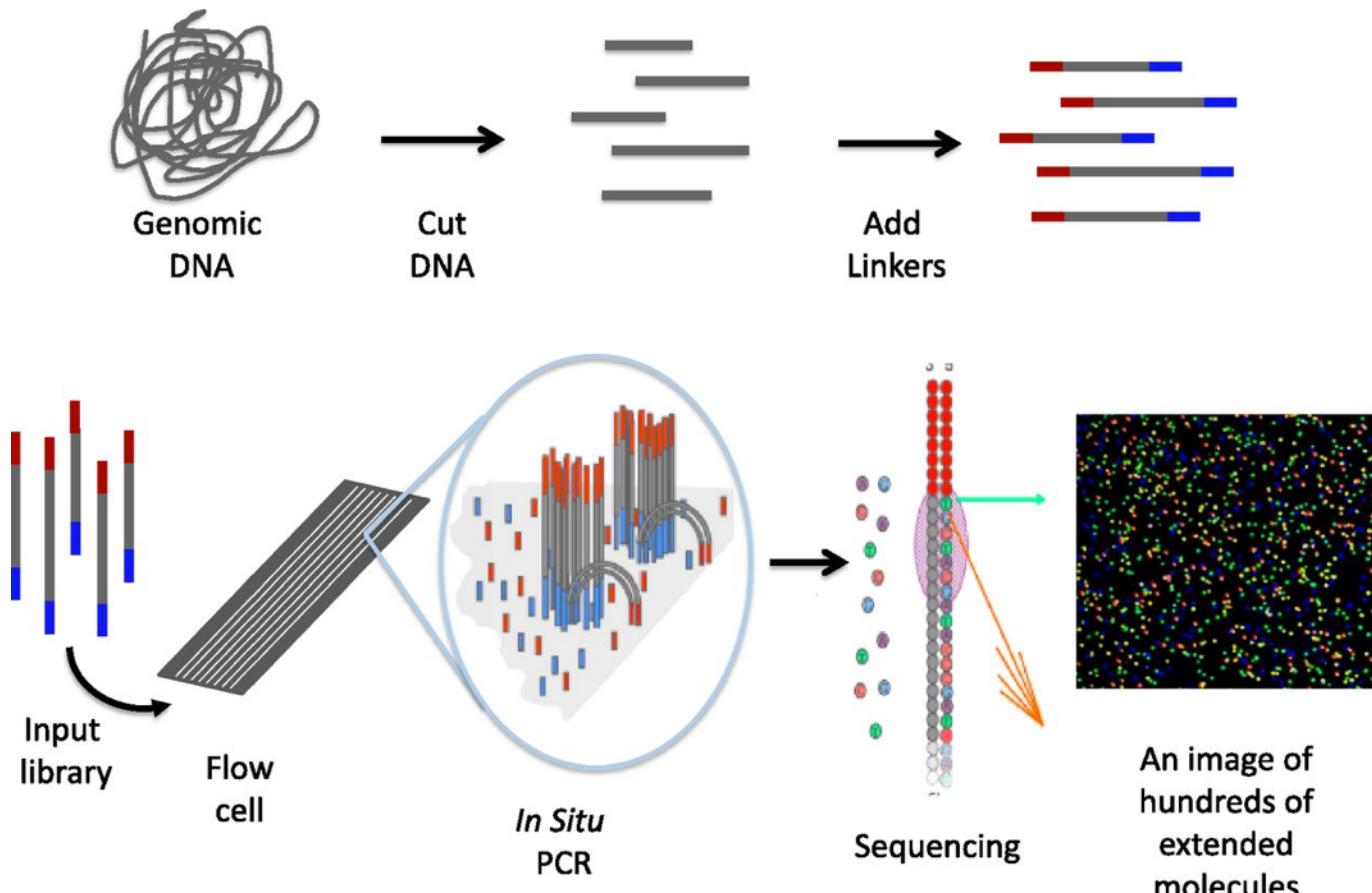
Editing level:
Fraction of edited transcripts

RNA editing in mammalian models



VGV mice = 100% Editing level at sites A-D

Second-generation sequencing



Applications of second-generation sequencing:

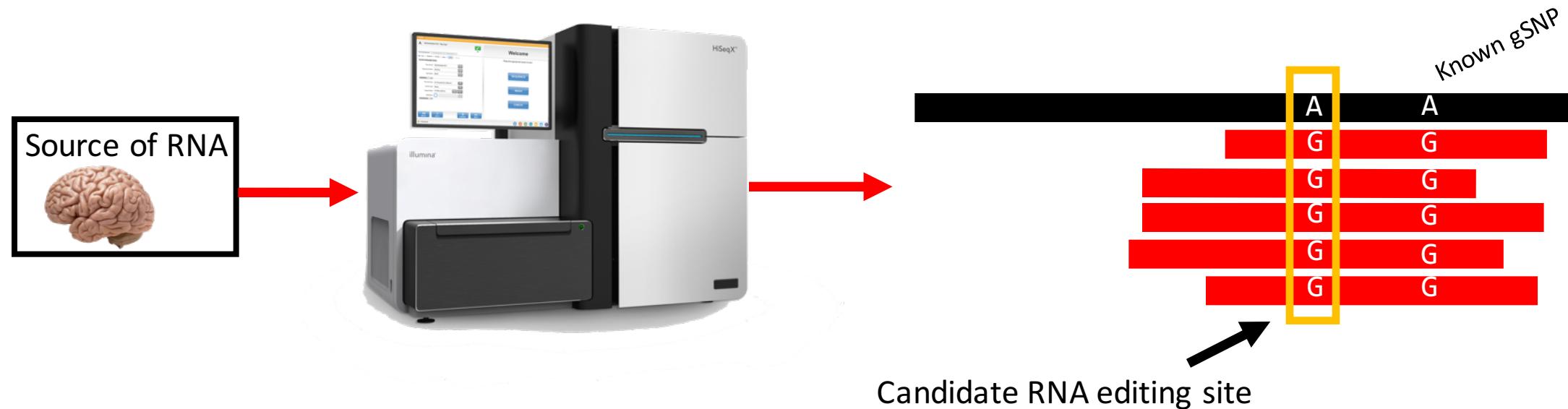
- Re-sequencing genomes to discover rare variants (segregating at low frequency; de-novo variants) *(re sequenced genome "leads")*
- Sequencing whole transcriptomes to quantify gene-expression levels

A-to-I RNA editing occurs at over a hundred million genomic sites, located in a majority of human genes

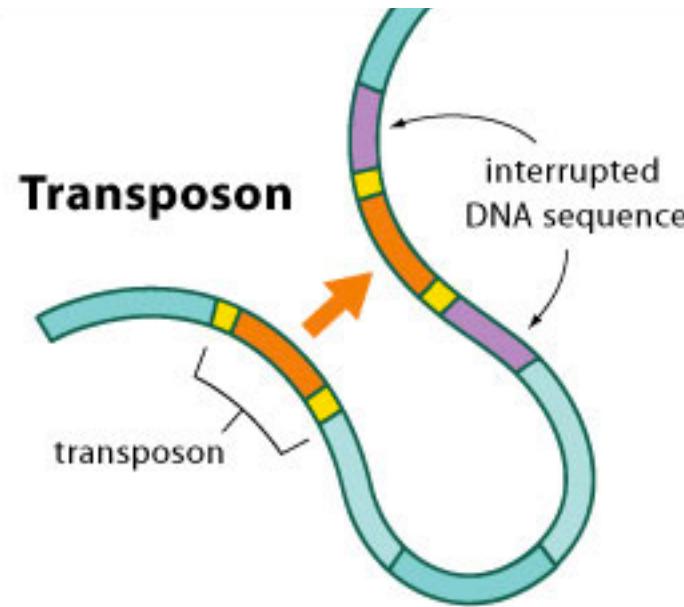
Lily Bazak,¹ Ami Haviv,¹ Michal Barak,¹ Jasmine Jacob-Hirsch,^{1,2} Patricia Deng,³ Rui Zhang,³ Farren J. Isaacs,⁴ Gideon Rechavi,^{2,5} Jin Billy Li,³ Eli Eisenberg,^{6,7,8} and Erez Y. Levanon^{1,7,8}

¹Mina and Everard Goodman Faculty of Life Sciences, Bar-Ilan University, Ramat Gan 52900, Israel; ²Cancer Research Center, Chaim Sheba Medical Center, Tel Hashomer 52621, Israel; ³Department of Genetics, Stanford University School of Medicine, Stanford, California 94305, USA; ⁴Department of Molecular, Cellular and Developmental Biology and Systems Biology Institute, Yale University, New Haven, Connecticut 06520, USA; ⁵Sackler School of Medicine, Tel Aviv University, Tel Aviv 69978, Israel; ⁶Raymond and Beverly Sackler School of Physics and Astronomy and Sagol School of Neuroscience, Tel Aviv University, Tel Aviv 69978, Israel

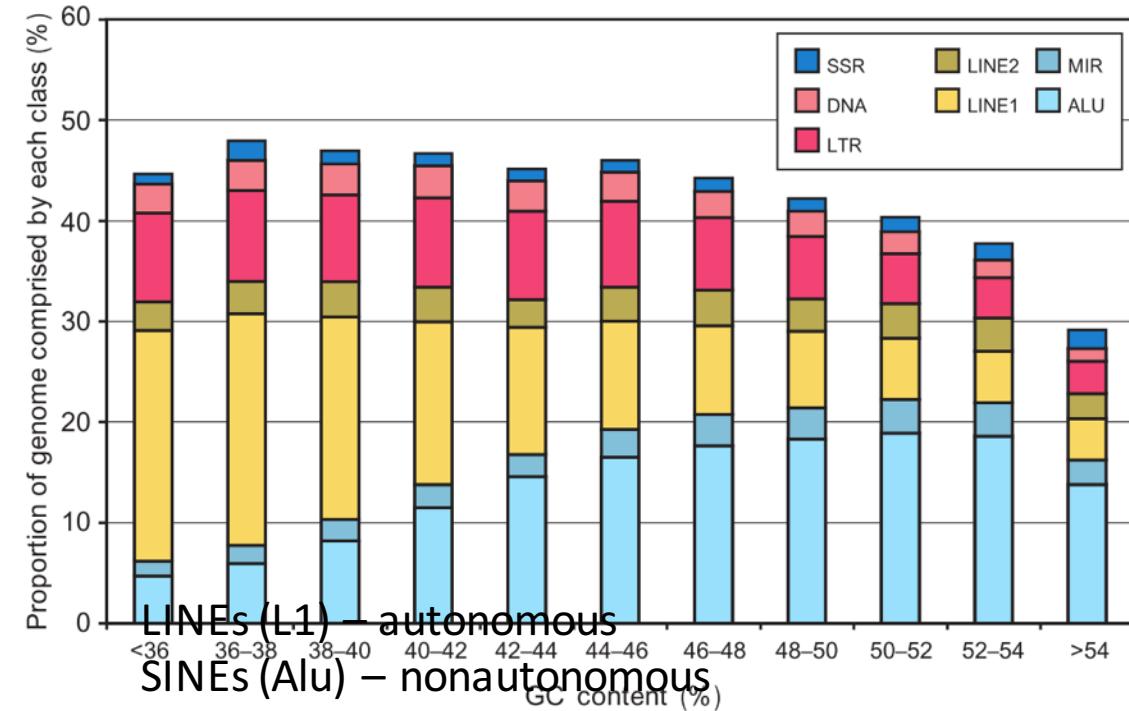
1.6 Million “edited sites”
>99% of all edited sites
in Alu elements
Nearly all adenosines
within Alu elements
appear editable.
Estimated 100M edited
sites.



Transposable elements (“Junk DNA”)



Human Genome

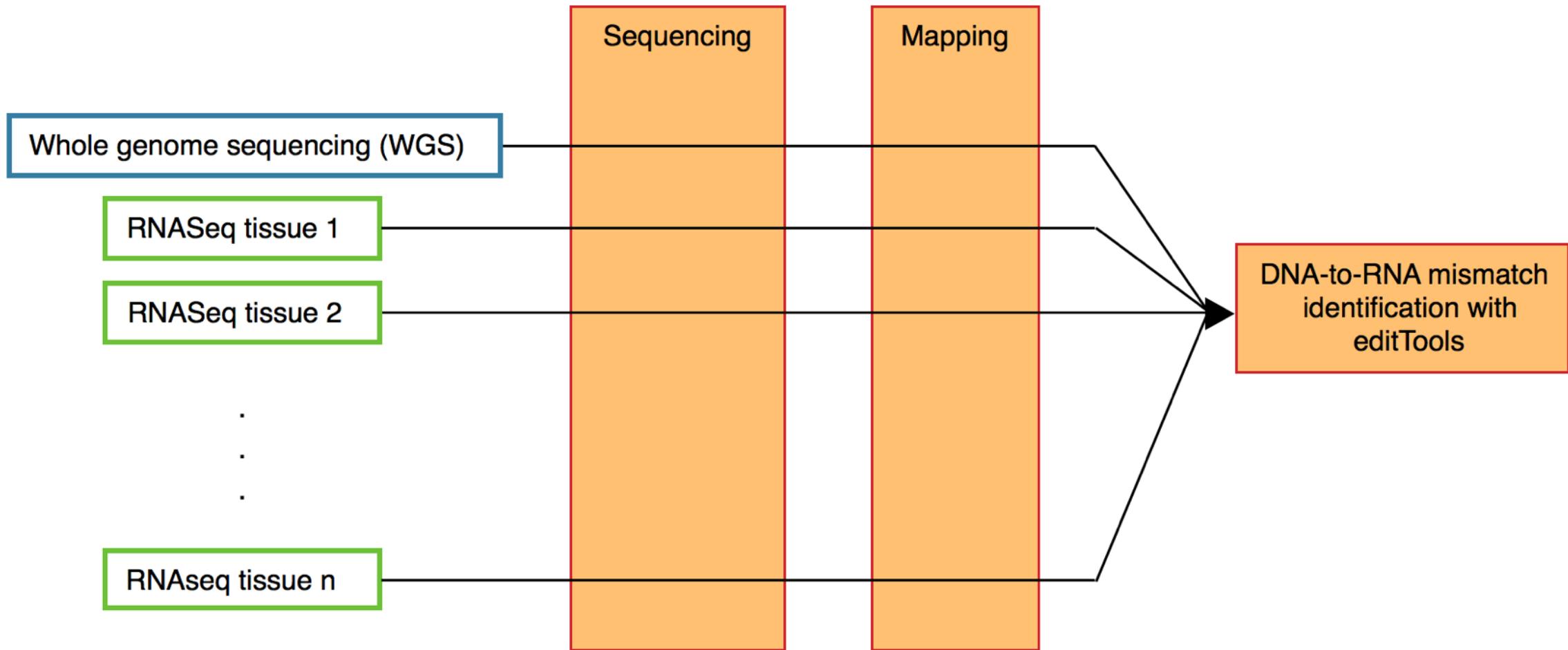


Currently...

- Most high-throughput RNA editing studies are done using human, other primates and mouse.



- Software designed to analyze RNA editing is nearly nonexistent



- Custom C++ libraries for processing of variant call format (VCF)
- *Rcpp: Seamless R and C++ Integration*

Variant call format

```

##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=.,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT Sample1 Sample2 Sample3
2 4370 rs6057 G A 29 . NS=2;DP=13;AF=0.5;DB;H2 GT:GQ:DP:HQ 0|0:48:1:52,51 1|0:48:8:51,51 1/1:43:5:..
2 7330 . T A 3 q10 NS=5;DP=12;AF=0.017 GT:GQ:DP:HQ 0|0:46:3:58,50 0|1:3:5:65,3 0/0:41:3
2 110696 rs6055 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27 2|1:2:0:18,2 2/2:35:4
2 130237 . T . 47 . NS=2;DP=16;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:56,51 0/0:61:2
2 134567 microsat1 GTCT G,GTACT 50 PASS NS=2;DP=9;AA=G GT:GQ:DP 0/1:35:4 0/2:17:2 1/1:40:3

```

```

class Rna {
    char genotype;
    int depth;
    int var_depth
    ...
}

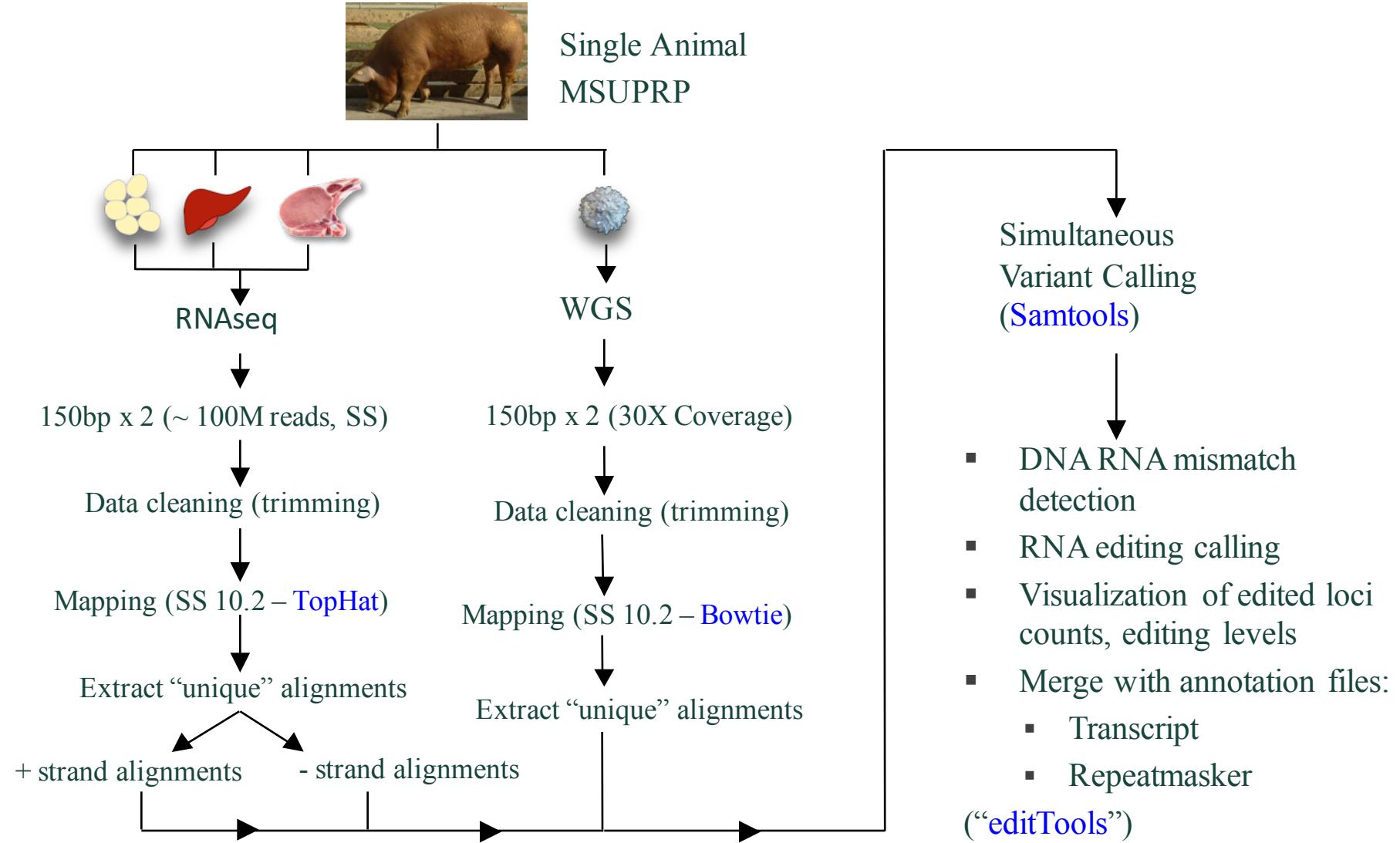
```

```

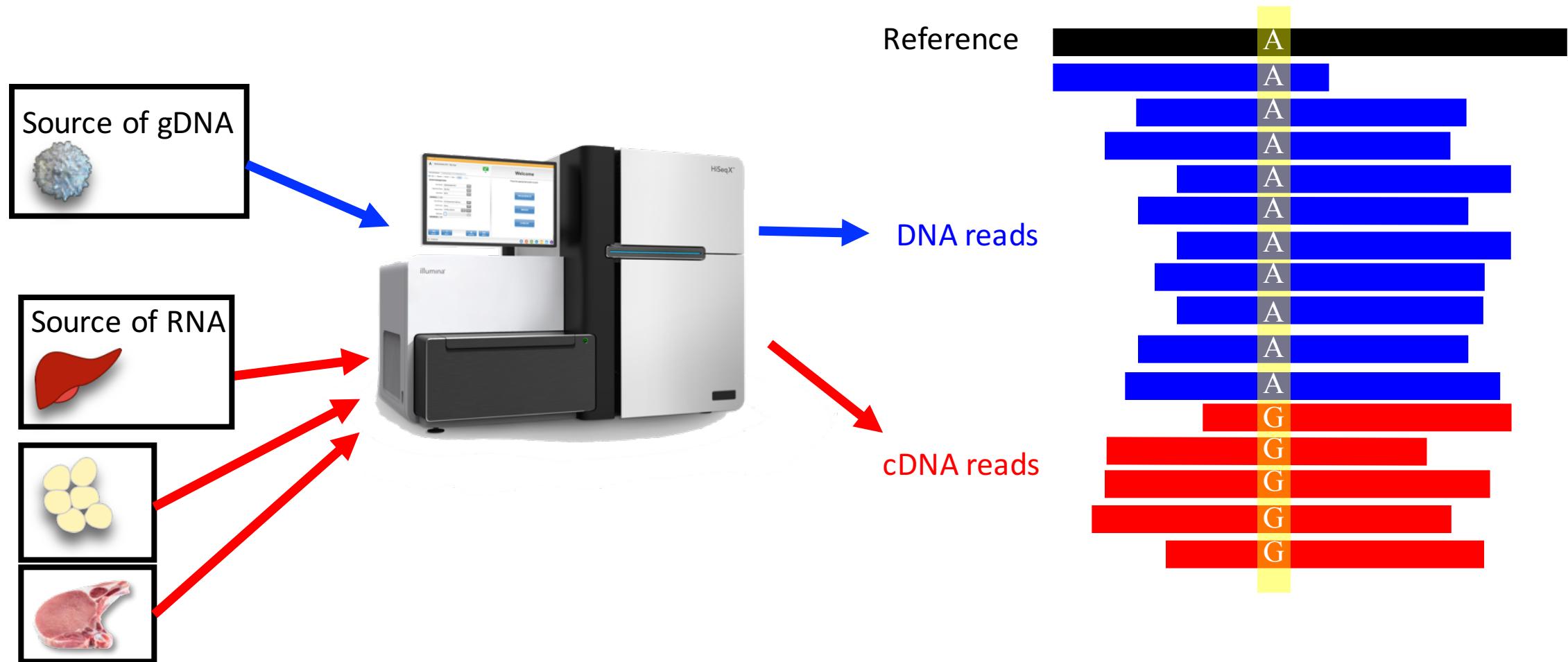
class Variant {
    char genotype;
    ...
    std::list<Rna> r_list;
}

```

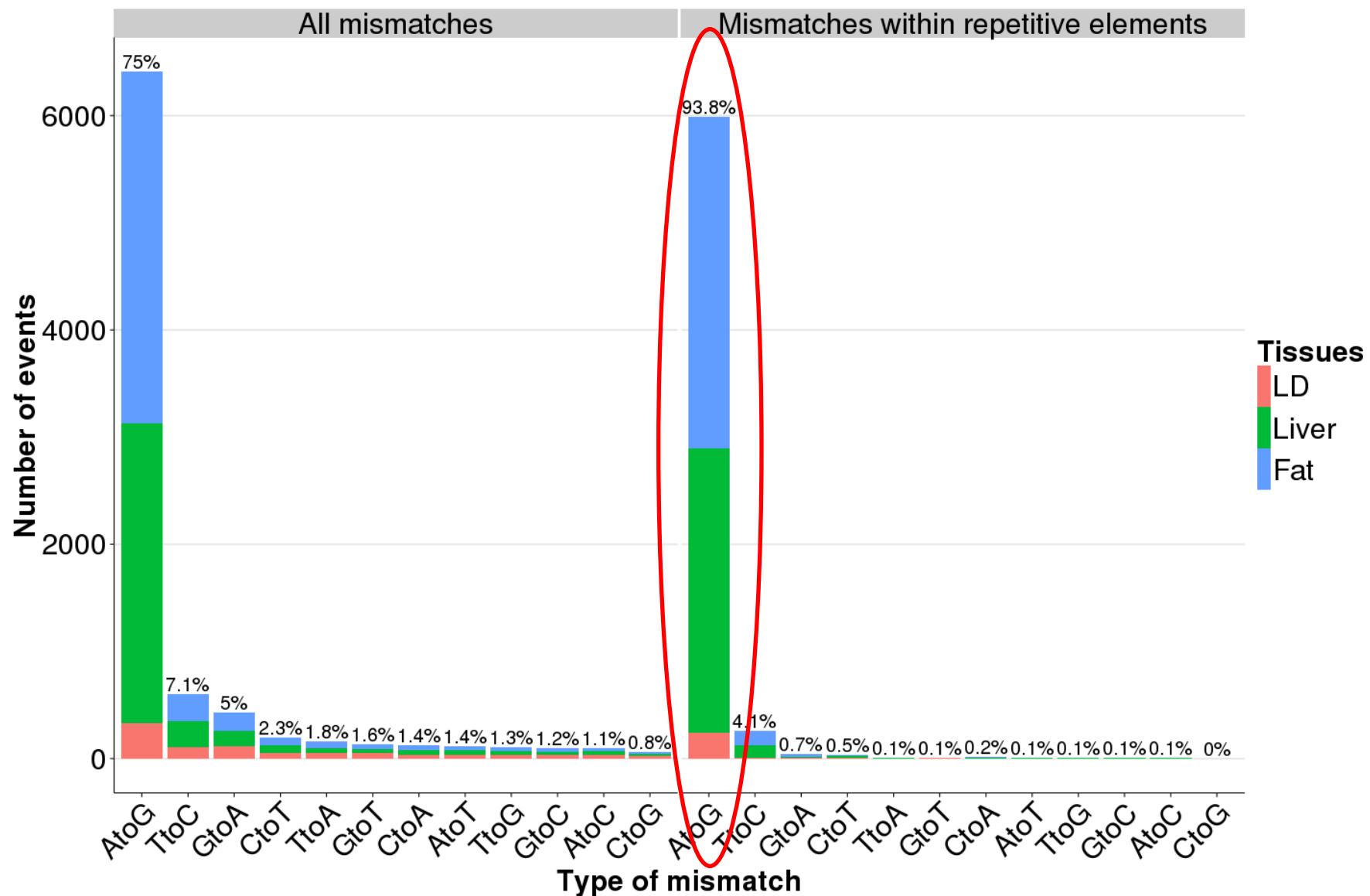
	DNA		RNA_1	
	FORMAT	Sample1	Sample2	Sample3
2	GT:GQ:DP:HQ	0 0:48:1:52,51	1 0:48:8:51,51	1/1:43:5:..
2	GT:GQ:DP:HQ	0 0:46:3:58,50	0 1:3:5:65,3	0/0:41:3
2	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
2	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:56,51	0/0:61:2
2	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3



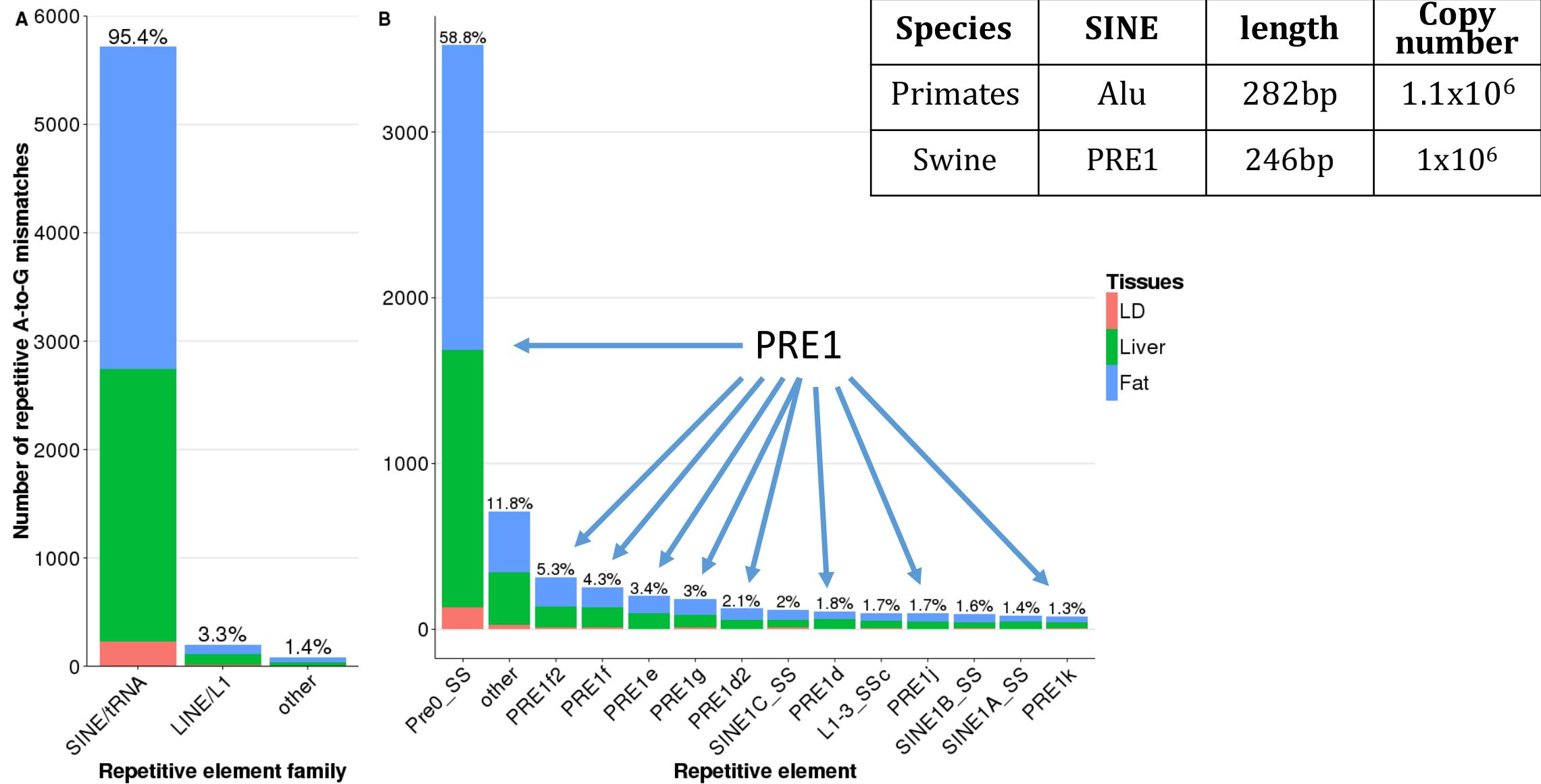
Calling RNA edited sites from WGS + RNASeq



Total DNA-to-RNA mismatch counts



Repetitive RNA editing distribution



Coding RNA editing distribution

- *GRIA2* – Normally edited at 100%. Lower editing levels associated with glioblastomas.
- *BLCAP* – Overediting associated with hepatocarcinoma.
- *NEIL1* – Unobservable in mouse

Position	Gene symbol/ID	AA	SIFT	Tissues
1:63408856	ENSSSCG00000029003	L/P	tolerated(1)	Fat LD Liver
1:125424444	ENSSSCG00000024660	Q/R	tolerated(1)	Fat LD Liver
11:22178068	<i>COG3</i>	I/V	tolerated(1)	Fat LD Liver
12:20231860	<i>AOC3</i>	Q/R	tolerated(1)	Liver
13:131377159	<i>EIF2B5</i>	Q/R	tolerated(1)	Fat
13:156760971	<i>UBE2B</i>	D/G	tolerated(0.48)	Fat LD Liver
13:206979572	<i>SON</i>	R/G	-	Fat
14:40832826	<i>PLBD2</i>	R/G	tolerated low confidence(0.12)	Fat
14:52398588	<i>IGLV-3</i>	E/G	tolerated(0.05)	Fat
14:59613334	<i>LYST</i>	S/G	-	LD
14:81796679	<i>OIT3</i>	S/G	tolerated(1)	Liver
15:59811585	<i>HNRNPA2B1</i>	L/P	tolerated(0.35)	Fat LD Liver
15:98217885	ENSSSCG00000028949	R/G	tolerated low confidence(1)	Fat LD Liver
16:29335640	ENSSSCG00000016869	N/D	tolerated(1)	Fat LD
16:42512978	<i>ELOVL7</i>	S/G	tolerated(1)	Fat
17:46041505	<i>BLCAP</i>	Y/C	deleterious(0)	Fat Liver
2:12622576	<i>LDHB</i>	I/M	tolerated(1)	Fat LD Liver
2:49316285	<i>ARNTL</i>	K/E	tolerated low confidence(1)	Liver
4:98044799	<i>COPA</i>	I/V	deleterious(0.02)	Fat
5:42375023	<i>KRR1</i>	I/T	deleterious(0.01)	Liver
6:92516721	<i>PTPRM</i>	K/R	tolerated(1)	Fat
6:146168578	<i>NDC1</i>	E/G	deleterious(0.01)	Liver
7:62951442	<i>NEIL1</i>	K/R	deleterious(0.02)	Fat LD
7:81602273	ENSSSCG00000002045	C/R	tolerated(1)	Fat LD Liver
7:102789222	<i>ACOT4</i>	T/A	tolerated(0.61)	Fat
7:129322238	<i>RPS21</i>	C/R	-	Fat LD Liver
8:28015971	ENSSSCG00000008767	H/R	tolerated(1)	Fat LD Liver
8:31629014	<i>TLR1</i>	I/V	tolerated(1)	Liver
8:32309809	<i>RPL9</i>	I/V	tolerated(0.4)	Fat
8:32309814	<i>RPL9</i>	E/G	deleterious(0.01)	Fat
8:48244993	<i>GRIA2</i>	Q/R	tolerated(0.07)	Fat
9:41146365	ENSSSCG00000023913	Q/R	deleterious(0.04)	Fat
9:74510703	ENSSSCG00000015294	K/R	tolerated(0.13)	Liver
9:83273454	<i>SLC25A13</i>	E/G	deleterious(0.02)	LD

Future goals

- *Evidence for substantial RNA editing transcriptome-wide among Sus scrofa PRE-1 SINE elements.* Intended for a “Report” in the journal RNA.
- Enhance the functionality and usability of *editTools*. Targeted for Bioconductor.
- FAANG – 6 Pigs. Each with WGS and RNASeq for *liver, fat, muscle, spleen, cortex, cerebellum, hypothalamus, and lung*.
- USDA-MARC – 3rd Generation Sequencing data. Compare to 2nd generation sequencing data from *hypothalamus, spleen, thymus and small intestine*.

Future of RNA editing

- Using RNA editing levels as predictors of quantitative traits
- $p \sim 1K$ or $1M$ or $100M$
- At the current cost of 2nd Gen sequencing, n would typically be very small.

	$j = 1$	$j = 2$	\dots	$j = p$
$i = 1$	1.0	0.0	\dots	
$i = 2$	0.85	\ddots		
\vdots	\vdots			
$i = n$				

