

## **Firehose AWG Runs**

**Purpose:** The Broad Institute GDAC will gladly coordinate with TCGA analysis working groups (AWGs) to provide custom Firehose runs tailored to their specific needs. This represents an evolution of Firehose, beyond its original mission of monthly runs intended for archival storage at the DCC and wide public consumption beyond TCGA, by providing indepth support to ongoing analysis efforts within TCGA. This provides several "real-time value-added benefits" to AWGs:

- <u>Currency</u>: pipelines can be run on the latest daily snapshot of data from the DCC, avoiding the time & sample lag of monthly runs
- <u>Flexibility</u>: additional runs can be easily performed on AWG-curated disease subtypes, and even include custom analyses
- <u>Speed</u>: custom AWG runs can be executed in only a few days time (excluding computationally intensive algorithms that may take >1 week to run) Familiarity: using the same internal Firehose machinery, and external-facing dashboards,
- Nozzle reports etc, already known to the community
- Easy Reproducibility: single system as starting point, everything kept, runs will soon be reference-able by DOI
- Scope: a stepping stone to open-access Firehose, that can be manipulated directly by TCGA researchers, instead of having runs curated only at the Broad

Our custom AWG data runs can also be used to define a baseline AWG data freeze. These freeze products are ideally suited for sharing data across the various centers participating in a given AWG. Furthermore, all of the output archives produced by our custom AWG runs are easily obtained with firehose\_get, in the same manner as the monthly runs.

		TCGA Firehose AWG analyses for THCA: 2013	03 18	
DiseaseType	AWG Run Dashboard	Maintained by <u>TCGA GDAC Team</u> (Broad Institute/MD Anderson Cancer Center/Harvard Medical School)		
GBM	<u>2013_02_17</u>	Unique Tumor Sample Counts	Data snapsh eading to free	
	2013_04_06	Tumor BCR Clinical CN LowP Methylation mRNA mRNAseq miR	miRseq RPPA N	
<u>LGG</u>	<u>2013_02_03</u>	THCA 500 318 432 94 500 0 409 0	411 224 3	
	<u>2013_01_16</u>	Download run results with <u>firehose get</u>		
HNSC	2013_03_30	Example download command: firehose_get awg_thca 2013_		
LUAD	<u>2013_02_07</u>	- Cancer Types	Analysis workfl On tumor p	
	<u>2012_11_15</u>		And 4 histol	
PANCAN8	<u>2012_08_25</u>	Table 1. Click "Browse" to view reports for a cancer type of interest. If you prefer to view reports on your containing all reports for a cancer type by clicking "Download".	own computer, you may downlo	
SKCM	<u>2013_01_16</u>	Cancer Type Cohort Reports HTML   Thyroid Adenocarcinoma THCA-TP 23 Browse		
	<u>2012_12_21</u>	firehose_getThyroid AdenocarcinomaTHCA-classical10BrowseOr browse fromThyroid AdenocarcinomaTHCA-follicular10Browse	Download	
STAD	<u>2013_04_17</u>	desktod	Download55 anDownloadII	
THCA	<u>2013_03_18</u> —			
	<u>2012_10_24</u>	http://gdac.broadinstitute.org/runs/awg_thca2013	03_18/	

# MAF Curation: Improving Consistency & Transparency via **Policy & Automation**

Historically, tracking MAFs has been difficult. There was no guarantee that a final AWG curated MAF would even make it to the DCC, let alone into Firehose. In order to keep the MAFs ingested into Firehose up-to-date, the maintainer had to track multiple AWG email lists, the AWGs' private wiki pages, the Jamboree site, and the DCC. Even then, it often came down to receiving an email from an analyst asking why the latest MAF wasn't available in Firehose.

Often these MAFs did not come from the DCC, requiring formatting cleanup and possible removal of germline data. Sometimes required data fields would be empty. Unvalidated MAFs and lack of a centralized repository led to inconsistent data.

To remedy this inconsistency, we now enforce the revised MAF data flow (pictured on the right). We now require that any MAF we ingest into Firehose MUST exist outside of it in a publically accessible fashion.

This enforcement has enabled us to better streamline our overall process. We have developed tools to automatically parse the MAFs that we mirror from the DCC and the MAFs we are ingesting into Firehose, summarize the information, and generate Confluence pages with these summaries laid out in an easy to read manner (partial examples on right).

### MAFs Available from the DCC as of 25 April 2013

	$\frac{1}{1}$										
	Tumor Type	Center	Tumor Samples	Normal Samples	Mutations	Columns	Archive Version	File Name			
'	BLCA	broad.mit.edu	136	145	44875	90	0.2.0	PR_TCGA_BLCA_PAIR_Capture_All_Pairs_QCPASS_v2.aggregated.capture.tcga.uuid.sc			
	BLCA	broad.mit.edu	28	28	7557	86	1.3.0	BLCA-28-original.aggregated.tcga.somatic.maf			
		genome.wustl.edu	107	109	7746	54	2.1.0	genome.wustl.edu_BRCA.IIIuminaGA_DNASeq.Level_2.2.0.0.maf			
	BRCA	genome.wustl.edu	510	546	29839	55	3.3.0	genome.wustl.edu_BRCA.IlluminaGA_DNASeq.Level_2.3.2.0.somatic.maf			
		genome.wustl.edu	776	801	47243	55	5.2.0	genome.wustl.edu_BRCA.IlluminaGA_DNASeq.Level_2.5.1.0.somatic.maf			
	CESC	broad.mit.edu	39	39	10020	90	1.4.0	PR_TCGA_CESC_PAIR_Capture_All_Pairs.aggregated.capture.tcga.uuid.somatic.			
	CESC	broad.mit.edu	40	40	11509	90	0.2.0	PR_TCGA_CESC_PAIR_Capture_All_Pairs_QCPASS_v2.aggregated.capture.tcga.uuid.sc			
	COAD	hgsc.bcm.edu	220	221	114594	47	1.4.0	hgsc.bcm.edu_COAD.IlluminaGA_DNASeq.1.somatic.maf			
	COAD	hgsc.bcm.edu	53	53	10986	41	1.5.0	hgsc.bcm.edu_COAD.SOLiD_DNASeq.1.somatic.maf			
	GBM	broad.mit.edu	291	291	22166	79	1.1.0	gbm_liftover.aggregated.capture.tcga.uuid.somatic.maf			

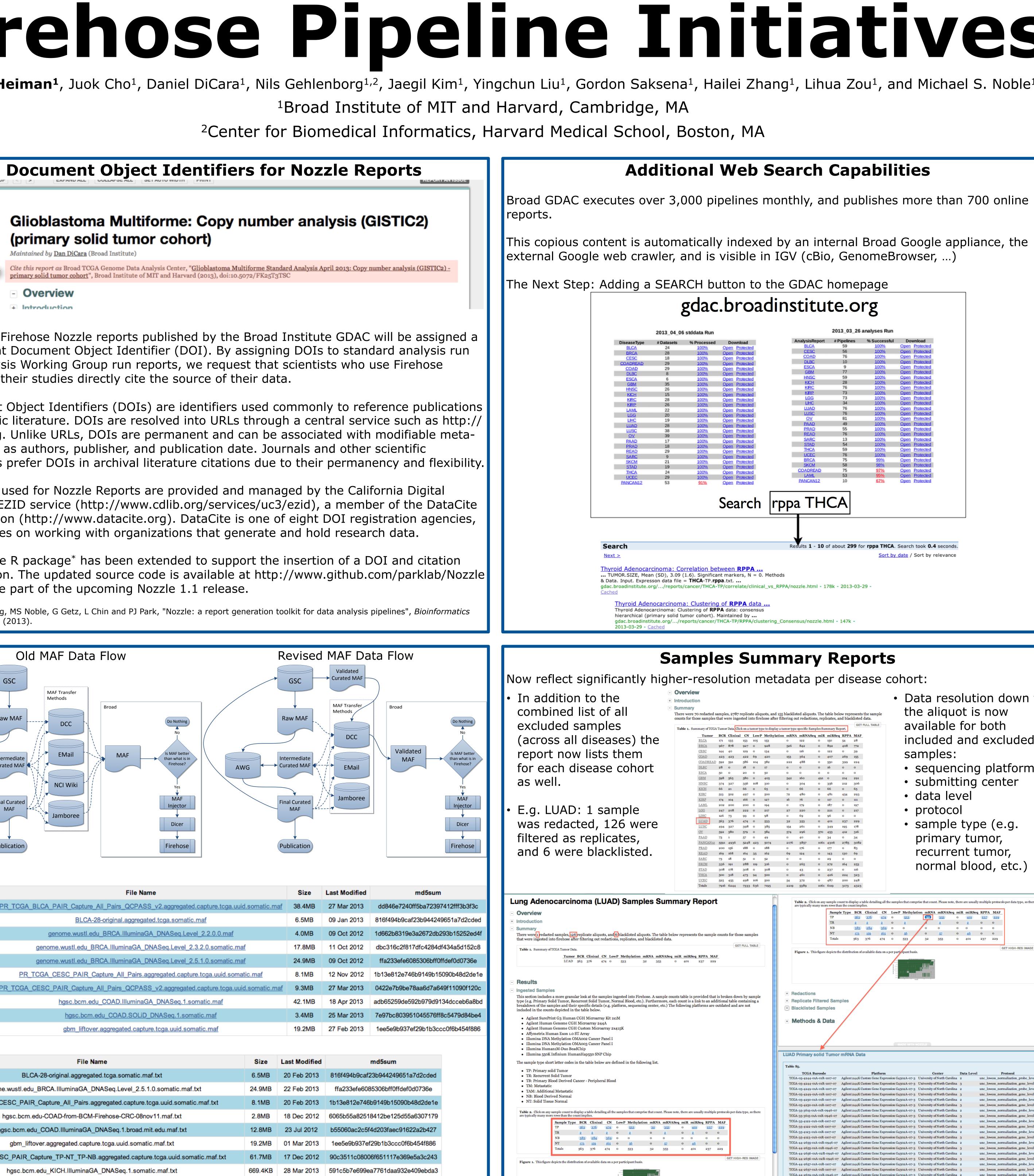
### MAFs Ingested into Broad GDAC Firehose as of 26 April 2013

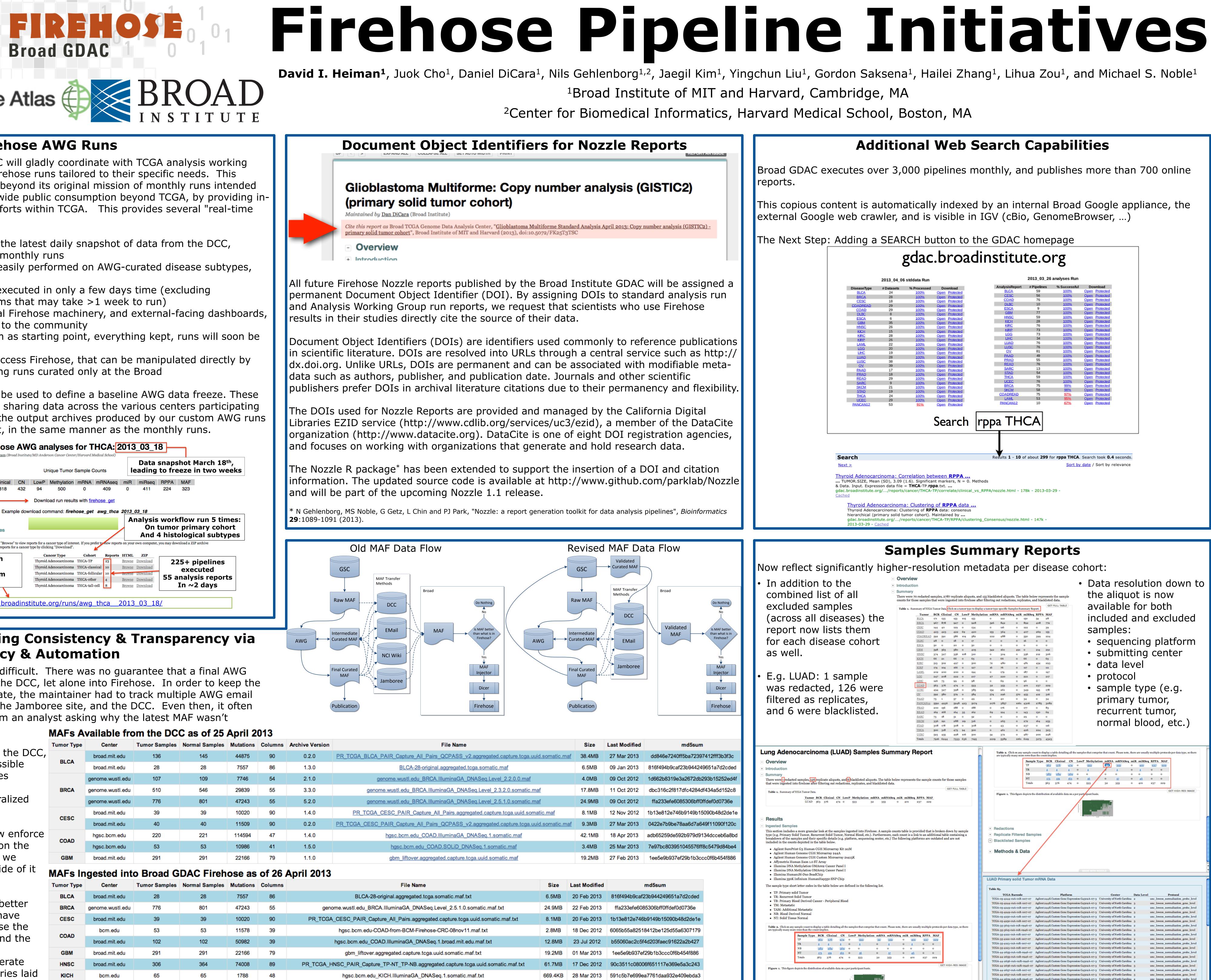
Tumor Type	Center	Tumor Samples	Normal Samples	Mutations	Columns	File Name	Size	Last Modified	md5sum
BLCA	broad.mit.edu	28	28	7557	86	BLCA-28-original.aggregated.tcga.somatic.maf.txt	6.5MB	20 Feb 2013	816f494b9caf23b944249651a7d2cded
BRCA	genome.wustl.edu	776	801	47243	55	genome.wustl.edu_BRCA.IlluminaGA_DNASeq.Level_2.5.1.0.somatic.maf.txt	24.9MB	22 Feb 2013	ffa233efe6085306bff0ffdef0d0736e
CESC	broad.mit.edu	39	39	10020	90	PR_TCGA_CESC_PAIR_Capture_All_Pairs.aggregated.capture.tcga.uuid.somatic.maf.txt	8.1MB	20 Feb 2013	1b13e812e746b9149b15090b48d2de1e
COAD	bcm.edu	53	53	11578	39	hgsc.bcm.edu-COAD-from-BCM-Firehose-CRC-08nov11.maf.txt	2.8MB	18 Dec 2012	6065b55a82518412be125d55a6307179
COAD	broad.mit.edu	102	102	50982	39	hgsc.bcm.edu_COAD.IIIuminaGA_DNASeq.1.broad.mit.edu.maf.txt	12.8MB	23 Jul 2012	b55060ac2c5f4d203faec91622a2b427
GBM	broad.mit.edu	291	291	22166	79	gbm_liftover.aggregated.capture.tcga.uuid.somatic.maf.txt	19.2MB	01 Mar 2013	1ee5e9b937ef29b1b3ccc0f6b454f886
HNSC	broad.mit.edu	306	364	74008	89	PR_TCGA_HNSC_PAIR_Capture_TP-NT_TP-NB.aggregated.capture.tcga.uuid.somatic.maf.txt	61.7MB	17 Dec 2012	90c3511c08006f651117e369e5a3c243
КІСН	bcm.edu	65	65	1788	48	hgsc.bcm.edu_KICH.IIIuminaGA_DNASeq.1.somatic.maf.txt	669.4KB	28 Mar 2013	591c5b7e699ea7761daa932e409ebda3
KIRC	bcm.edu	176	176	8768	32	BCM-BI-Renal-v1.3.hgsc.bcm.edu.maf.txt	1.8MB	23 Nov 2011	ed030b92cf604abfb3fb4a7d9186d235
NING	broad.mit.edu	297	297	95839	84	BI_and_BCM_1.4.aggregated.tcga.broad.mit.edu.maf.txt	58.9MB	30 Jul 2012	2ab7c89d11a077618050fb4e6f5390ca

hot March 18<sup>th</sup>, eze in two weeks

flow run 5 times: primary cohort ological subtypes ad a ZIP archive

25+ pipelines executed nalysis reports In ~2 days





 Data resolution down to the aliquot is now available for both included and excluded sequencing platform submitting center data level protocol • sample type (e.g. primary tumor recurrent tumor, normal blood, etc.) Table 2. Click on any sample count to display a table detailing all the samples that comprise that count. Please note, there are usually multiple protocols per data type, so ther Sample Type BCR Clinical CN LowP Methylation mRNA mRNASeq miR miRSeq RPPA MAF 353 0 401 237 229 TR 2 2 2 0 2 0 2 0 0 0 0 0 0 0 0 NT <u>173</u> <u>159</u> <u>165</u> 0 <u>56</u> 0 <u>57</u> 0 <u>46</u> 0 0 563 376 474 0 533 32 353 0 401 237 229 GET HIGH-RES IMAGE Data Level unc\_lowess\_normalization\_probe\_leve

TCGA-44-2659-01A-01R-0946-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 2 TCGA-44-2659-01A-01R-0946-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 3 TCGA-44-2661-01A-01R-1107-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 2

unc\_lowess\_normalization\_probe\_level unc lowess normalization gene level unc lowess normalization probe leve unc\_lowess\_normalization\_gene\_leve unc\_lowess\_normalization\_probe\_leve unc\_lowess\_normalization\_gene\_level unc\_lowess\_normalization\_probe\_lev unc\_lowess\_normalization\_gene\_leve unc\_lowess\_normalization\_probe\_leve unc\_lowess\_normalization\_gene\_lev unc\_lowess\_normalization\_probe\_lev unc\_lowess\_normalization\_gene\_leve unc\_lowess\_normalization\_probe\_leve unc\_lowess\_normalization\_gene\_level unc\_lowess\_normalization\_probe\_level unc\_lowess\_normalization\_gene\_level unc\_lowess\_normalization\_probe\_leve unc\_lowess\_normalization\_gene\_level unc\_lowess\_normalization\_probe\_level TCGA-44-2661-01A-01R-1107-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 3 unc\_lowess\_normalization\_gene\_level TCGA-44-2662-01A-01R-0946-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 2 unc\_lowess\_normalization\_probe\_level TCGA-44-2662-01A-01R-0946-07 Agilent 244K Custom Gene Expression G4502A-07-3 University of North Carolina 3 unc\_lowess\_normalization\_gene\_level