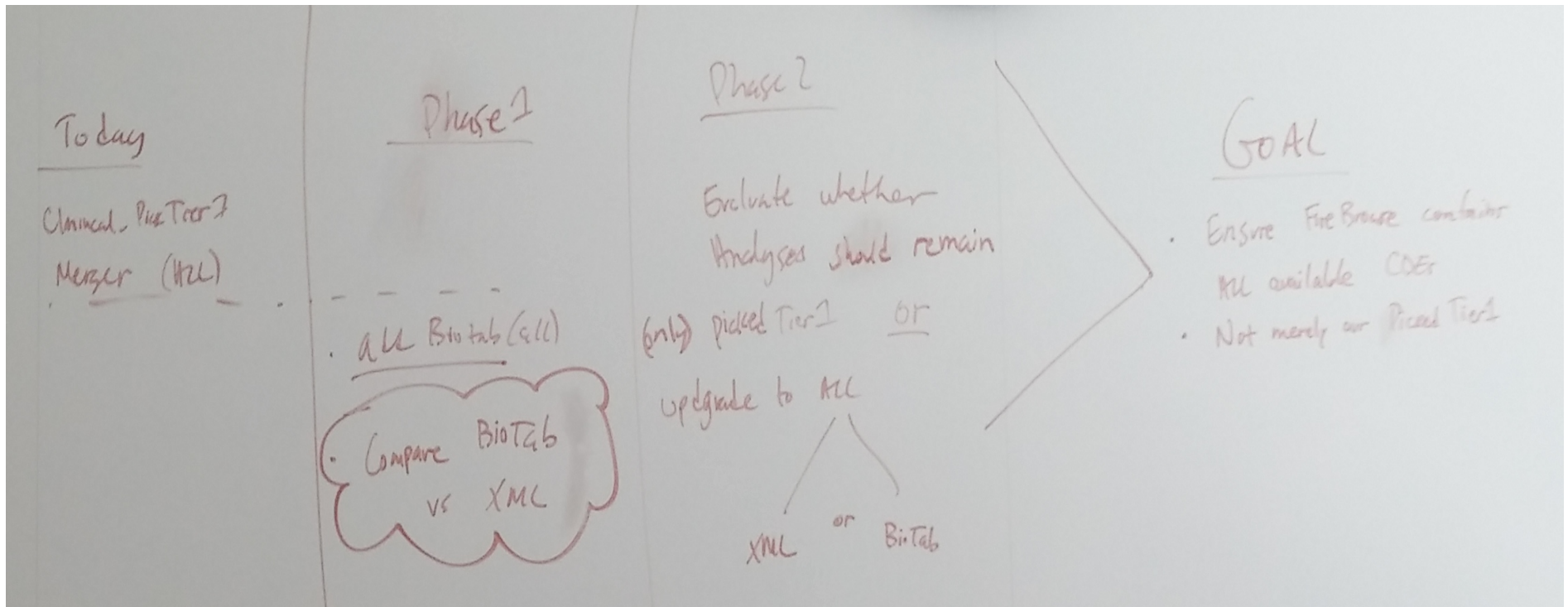


Clinical data meeting – July 21 2015



TO DO

- Phase 1. Test module for BioTab(all) will go parallel with existing modules using XML
-> will compare BioTab(all) vs XML(all) in the merger output.
- Phase 2. will evaluate if remain on picked or upgrade to all(of XML/BioTab)

Unique parameter count table of XML(25) and BioTab(34) for 25 common diseases

- XML mostly has more parameters than BioTab because multiple versions are in slightly different param names in xml data.
- It is hard to obtain exactly unique parameter set to compare in xml data and the unique sample count in xml is rough number.

: BioTab clinical format has two headers.

Clinical

disease	format	xml.uniq.p.coun	biotab.uniq.p.count_h1.c	biotab.p.count_h2.c
ACC	clinical	170	164	162
BLCA	clinical	171	169	160
BRCA	clinical	204	203	194
CESC	clinical	230	244	234
CHOL	clinical	165	175	162
COAD	clinical	154	166	159
COADREAD	clinical	155	NA	NA
DLBC	clinical	155	158	184
ESCA	clinical	163	169	160
FPPP	clinical	301	NA	NA
GBM	clinical	97	115	111
GBMLGG	clinical	150	NA	NA
HNSC	clinical	164	171	161
KICH	clinical	141	136	128
KIPAN	clinical	146	NA	NA
KIRC	clinical	141	147	138
KIRP	clinical	141	150	139
LAML	clinical	166	77	78
LGG	clinical	143	161	156
LIHC	clinical	168	172	159
LUAD	clinical	151	198	188
LUSC	clinical	152	196	186
MESO	clinical	147	154	146
OV	clinical	144	125	122
PAAD	clinical	159	180	173

Auxiliary

disease	format	xml.uniq.p.coun	biotab.uniq.p.coun
CESC	auxiliary	16	12
COAD	auxiliary	15	12
COADREAD	auxiliary	15	NA
ESCA	auxiliary	15	12
HNSC	auxiliary	16	12
PAAD	auxiliary	15	12

Biospecimen

: biospecimen parameters are not proper for correlation analysis and will be excluded.

disease	format	xml.uniq.p.coun	biotab.uniq.p.count
ACC	biospecimen	160	127
BLCA	biospecimen	168	127
BRCA	biospecimen	147	127
CESC	biospecimen	160	127
CHOL	biospecimen	166	127
COAD	biospecimen	158	127
COADREAD	biospecimen	144	NA
DLBC	biospecimen	155	127
ESCA	biospecimen	160	127
FPPP	biospecimen	137	106
GBM	biospecimen	165	127
GBMLGG	biospecimen	173	NA
HNSC	biospecimen	160	127
KICH	biospecimen	158	127
KIPAN	biospecimen	166	NA
KIRC	biospecimen	158	127
KIRP	biospecimen	166	127
LAML	biospecimen	166	83
LGG	biospecimen	171	127
LIHC	biospecimen	166	127
LUAD	biospecimen	165	127
LUSC	biospecimen	164	127
MESO	biospecimen	157	127
OV	biospecimen	160	127
PAAD	biospecimen	164	127

BioTab Param groups table for all diseases

- BioTab data files saved in .txt by different parameter groups.
- Param groups in yellow color are mostly useful for correlation analysis.
- XML data parameters are not grouped and mixed up so providing all CDEs leads meaningless parameters to be used in correlation analysis.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ									
1		acc	blca	blnp	blp	brca	cesc	chol	cntl	coad	dlbc	esca	fppp	gbm	hnsc	kich	kirc	kirp	laml	lcll	lcmi	lgg	lihc	lndlbc	lnnh	luad	lusc	meso	misc	ov	paad	pcpg	prad	read	sarc	skcm	stad	tgct	thca	thym	ucec	ucs	uvm									
2	aliquot	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1						
3	analyte	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1				
4	cqcf	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1				
5	diagnostic_slides	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
6	normal_control	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
7	portion	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
8	protocol	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
9	sample	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1			
10	shipment_portion	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
11	slide	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
12	tumor_sample	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
13	drug	1	1	0	0	1	1	1	0	1	0	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
14	nte	1	1	0	0	1	1	1	0	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
15	omf_v4.0	1	1	0	0	1	1	1	0	1	0	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
16	patient	1	1	0	0	1	1	1	0	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
17	radiation	1	1	0	0	1	1	1	0	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	
18	follow_up	1	1	0	0	1	1	1	0	1	1	1	0	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
19	control	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Param group

Biotab data usage in CHOL AWG

: They use several param groups only for analysis.

Files » TCGA live » CHOL » bio

Synapse ID: syn2671816

Conditions
Param group

Name	Conditions	Param group
nationwidechildrens.org_CHOL_bio.aliquot.tsv		
nationwidechildrens.org_CHOL_bio.analyte.tsv		
nationwidechildrens.org_CHOL_bio.drug.tsv		
nationwidechildrens.org_CHOL_bio.followup.tsv		
nationwidechildrens.org_CHOL_bio.patient.tsv		
nationwidechildrens.org_CHOL_bio.portion.tsv		
nationwidechildrens.org_CHOL_bio.radiation.tsv		
nationwidechildrens.org_CHOL_bio.sample.tsv		

nationwidechildrens.org_biospecimen_aliquot_chol.txt
nationwidechildrens.org_biospecimen_analyte_chol.txt
nationwidechildrens.org_biospecimen_cqcf_chol.txt
nationwidechildrens.org_biospecimen_diagnostic_slides_chol.txt
nationwidechildrens.org_biospecimen_normal_control_chol.txt
nationwidechildrens.org_biospecimen_portion_chol.txt
nationwidechildrens.org_biospecimen_protocol_chol.txt
nationwidechildrens.org_biospecimen_sample_chol.txt
nationwidechildrens.org_biospecimen_shipment_portion_chol.txt
nationwidechildrens.org_biospecimen_slide_chol.txt
nationwidechildrens.org_biospecimen_tumor_sample_chol.txt
nationwidechildrens.org_clinical_cqcf_chol.txt
nationwidechildrens.org_clinical_drug_chol.txt
nationwidechildrens.org_clinical_follow_up_v4.0_chol.txt
nationwidechildrens.org_clinical_follow_up_v4.0_nte_chol.txt
nationwidechildrens.org_clinical_nte_chol.txt
nationwidechildrens.org_clinical_omf_v4.0_chol.txt
nationwidechildrens.org_clinical_patient_chol.txt
nationwidechildrens.org_clinical_radiation_chol.txt

Not used?

: For the picker pipeline,

- only parameters in param groups of drug, follow_up, patient, radiation are proper for correlation analysis.
- Compared to BioTab data, xml data doesn't categorize the parameters by the groups and are mixed up.

Params in param groups of sample, portion, aliquot, analyte: not proper for correlation analysis

Sample

bcr_patient_uuid	bcr_sample_barcode	bcr_sample_uuid	NCNNCT_0thMethONSP	b
iospecimen_sequence	composition	current_weight	days_to_collection	days_t
o_sample_procurement	freezing_method	initial_weight	intermediate_dimension	is_ffp
e	longest_dimension	method_of_sample_procurement	oct_embedded	other_
method_of_sample_procurement	pathology_report_file_name	pathology_report_uuid		
reservation_method	sample_type	sample_type_id	shortest_dimension	time_b
etween_clamping_and_freezing	time_between_excision_and_freezing	tissue_type		t
tumor_descriptor	vial_number			

portion

bcr_patient_uuid	bcr_sample_barcode	bcr_portion_barcode	bcr_portion_uu		
id	LCE	date_of_creation	is_ffpe	portion_number	portion_sequence
weight					

aliquot

bcr_patient_uuid	bcr_sample_barcode	bcr_analyte_barcode	bcr_analyte_uu			
id	a260_a280_ratio	amount	analyte_type	analyte_type_id	concentration	gel_im
age_file	is_derived_from_ffpe	normal_tumor_genotype_match	pcr_amplificat			
ion_successful	ratio_28s_18s	rvalue	spectrophotometer_method	subpor		
tion_sequence	well_number					

analyte

bcr_patient_uuid	bcr_sample_barcode	bcr_aliquot_barcode	bcr_aliquot_uu		
id	amount	biospecimen_barcode_bottom	center_id	concentration	date_o
f_shipment	is_derived_from_ffpe	plate_column	plate_id	plate_row	q
quantity	source_center	volume			

OPS meeting agenda for clinical update - 2015

June 15, 2015

1. Welcome Tim
2. iCoMut for RAS pathway: status?
 - a. check integration—at present there are aggregate failures
3. iCoMut for SARC run: presentation deferred until next week, but still have legwork to do
4. Review & parcel out unanswered GDAC list questions
5. Review CompBio projects:
 - a. [Update clinical pipelines](#) to potentially use BioTAB, and reveal ALL **clinical** parameters--not just picked.
 - b. Recent / Upcoming AWG runs:
 - i. PAAD:
 - ii. CHOL:
 - c. Bayesian NMF
 - d. MGH: horizon?
 - e. More?
6. Revisit [prioritizing next 6-12 month dev window](#)

f. Where would Sam fit?

5. [viewGene widnet RFEs](#)

June 1, 2015

lisease cohorts with no RSEM values

6. iCoMut for RAS pathway
7. **Clinical** data:
 - a. why continue to use XML over biotab?
 - b. Per Gordon: XML was preferred over biotab because
 - i. it's the source data
 - ii. from which biotab was generated (by DCC)
 - iii. and biotab was more buggy
 - iv. and it comes in packages of 5-6 files per cohort, which was harder to merge
 - c. we need to provide ALL of **clinical**, not just our picked subset

Clinical data update agenda - 2015

- All clinical features
 - current xml data
 - BioTab data
 - ⇒ Will try BioTab first and then will compare with xml.
- Up-to-date data
 - There are more parameters in followup data versioned (see slide 7)
 - ⇒ Will update pipeline, SelectionFileGenerator for all followup parameters such as radiation_therapy.
- BioTab data
 - XML: 3 txt files generated from the dicer running XML parser
 - BioTab: Can BioTab also be generated from the dicer? (Not yet)
 - ⇒ Will check if BioTab can be added in a parameter in Clinical_Merger.
- MSI data
 - According to Gordon's email (see next slide 4), there are formats unprocessed having MSI data.
 - ⇒ Will skipped them as our dicer doesn't work for them.

> Gordon recently reported that he found new data type having MSI info and the data type was not processed due to ingestor failures during the 2014_10_06 GDAC ops meeting.

> According to Gordon, the new data types are as below.

> The files live here:

> 3:24pm gsaksena@voncotator2-dev /xchip/gdac_data/dcc_mirror3/platform_link \$ ls -1 *micro*

> coad__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__anonymous@

> coad__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__tcga4yeo@

> read__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__anonymous@

> read__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__tcga4yeo@

> stad__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__anonymous@

> stad__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__tcga4yeo@

> ucec__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__anonymous@

> ucec__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__tcga4yeo@

> ucs__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__anonymous@

> ucs__cgcc__nationwidechildrens.org__microsat_i__fragment_analysis__tcga4yeo@

> The corresponding unprocessed flags from the dicer are here, which indicates the ingestor is ignoring this data type:

> 3:27pm gsaksena@cga-cdkn2a /xchip/gdac_data/normalized/sdrf_cache \$ ls -1

*unprocessed

> nationwidechildrens.org_COAD.microsat_i.mage-tab.1.6.0.txt.unprocessed

> nationwidechildrens.org_READ.microsat_i.mage-tab.1.7.0.txt.unprocessed

> nationwidechildrens.org_STAD.microsat_i.mage-tab.1.0.0.txt.unprocessed

> nationwidechildrens.org_UCEC.microsat_i.mage-tab.1.6.0.txt.unprocessed

> nationwidechildrens.org_UCS.microsat_i.mage-tab.1.0.0.txt.unprocessed

Parameters

Parameter Name	Parameter Expression	Mode	Default Value
sdrfName	<input type="text"/>	Simple_Expression <input type="button" value="v"/>	<input type="text"/>
doNotCreateManifest	<input type="text" value="-m"/>	Literal <input type="button" value="v"/>	<input type="text"/>
samplestamp*	<input type="text" value="samplestamp"/>	Simple_Expression <input type="button" value="v"/>	<input type="text"/>
annotationids*	<input type="text" value="clin__bio__nationwidechilk"/>	Literal <input type="button" value="v"/>	<input type="text"/>
datapathsfile1*	<input clin__bi"="" type="text" value="samples.choose(["/>	Complex_Expression <input type="button" value="v"/>	<input type="text"/>
inFile_clinical_clin*	<input clin__bi"="" type="text" value="samples.choose(["/>	Complex_Expression <input type="button" value="v"/>	<input type="text"/>
inFile_biospecimen_clin	<input clin__bi"="" type="text" value="samples.choose(["/>	Complex_Expression <input type="button" value="v"/>	<input type="text"/>
inFile_auxiliary_clin	<input clin__bi"="" type="text" value="samples.choose(["/>	Complex_Expression <input type="button" value="v"/>	<input type="text"/>
outPrefix*	<input type="text" value="sample_set_id"/>	Simple_Expression <input type="button" value="v"/>	<input type="text"/>

Additional parameters?

Vital_status parameters in the PAAD-TP.clin.merged.txt

	A	B	C	D	E	F	G
1	patient.bcr_patient_barcode	tcga-2j-aab1	tcga-2j-aab4	tcga-2j-aab6	tcga-2j-aab8	tcga-2j-aab9	tcga-2j-aaba
2	patient.vital_status	dead	alive	dead	alive	dead	dead
3	patient.follow_ups.follow_up.vital_status	NA	alive	NA	alive	NA	dead
4	patient.follow_ups.follow_up-2.vital_status	NA	NA	NA	NA	NA	NA
5	patient.follow_ups.follow_up-3.vital_status	NA	NA	NA	NA	NA	NA
Final values for vital_status		Dead	Alive	Dead	Alive	Dead	dead
Where the value comes from		Primary param	Followup version1	Primary param	Followup version1	Primary param	Followup version1

- As seen in the table above,
 1. Using the latest followup version 3 is not correct
 2. Using only one of versions is not correct
 3. NA in the latest version followup data means there is no update from the previous version.
 4. Thus, data values should be the combination of all versions.
 5. Also, It is not helpful to use the DCC upload date parameter to pick up the latest followup because the values are NA if there is no change from the previous version.

=> For the AWG run generator, recommend to check if there are new version of followup data ingested in the .clin.merged.txt

Radiation_therapy parameters in the COADREAD-TP.clin.merged.txt

1	V1
565	patient.follow_ups.follow_up-2.new_tumor_events.new_tumor_event.additional_radiation_therapy
576	patient.follow_ups.follow_up-2.radiation_therapy
589	patient.follow_ups.follow_up-3.new_tumor_events.new_tumor_event.additional_radiation_therapy
600	patient.follow_ups.follow_up-3.radiation_therapy
613	patient.follow_ups.follow_up-4.new_tumor_events.new_tumor_event.additional_radiation_therapy
624	patient.follow_ups.follow_up-4.radiation_therapy
637	patient.follow_ups.follow_up.new_tumor_events.new_tumor_event-2.additional_radiation_therapy
645	patient.follow_ups.follow_up.new_tumor_events.new_tumor_event.additional_radiation_therapy
656	patient.follow_ups.follow_up.radiation_therapy
681	patient.new_tumor_events.new_tumor_event.additional_radiation_therapy
704	patient.radiation_therapy
711	patient.radiations.radiation-2.days_to_radiation_therapy_end
712	patient.radiations.radiation-2.days_to_radiation_therapy_start
729	patient.radiations.radiation-3.days_to_radiation_therapy_end
730	patient.radiations.radiation-3.days_to_radiation_therapy_start
747	patient.radiations.radiation-4.days_to_radiation_therapy_end
748	patient.radiations.radiation-4.days_to_radiation_therapy_start
765	patient.radiations.radiation.days_to_radiation_therapy_end
766	patient.radiations.radiation.days_to_radiation_therapy_start

- Here, followup version 4 is the latest version but the version number is not chronological order. Our SelectionFileGenerator will check all values of each parameter and assign the latest value.

For the variable, "additional_phamaceutical_therapy"

*** CASE1. How to order ?? ** #

```
patient.follow_ups.follow_up-2.new_tumor_events.new_tumor_event.additional_pharmaceutical_therapy
patient.follow_ups.follow_up-3.new_tumor_events.new_tumor_event.additional_pharmaceutical_therapy
patient.follow_ups.follow_up-4.new_tumor_events.new_tumor_event.additional_pharmaceutical_therapy
patient.follow_ups.follow_up.new_tumor_events.new_tumor_event-2.additional_pharmaceutical_therapy
patient.follow_ups.follow_up.new_tumor_events.new_tumor_event.additional_pharmaceutical_therapy
patient.new_tumor_events.new_tumor_event.additional_pharmaceutical_therapy
```

Is there follow_up nodes? yes

-> find the follow_up nodes

1. find a node not matched with "follow_up"

2. find matches with ".follow_up."

* Is there multiple matches?

-> find the upperNode

2.1. find matches with ".upperNode."

2.2. find matches with ".upperNode-#."

3. find matches with ".follow_up-#." |

For the variable, "days_to_radiation_therapy_start"

*** CASE2. How to order ?? ** #

```
patient.radiations.radiation-2.days_to_radiation_therapy_start
```

```
patient.radiations.radiation-3.days_to_radiation_therapy_start
```

```
patient.radiations.radiation-4.days_to_radiation_therapy_start
```

```
patient.radiations.radiation.days_to_radiation_therapy_start
```

Is there follow_up nodes? no

-> find the upperNode

1. find matches with ".upperNode."

2. find matches with ".upperNode-#."

After setting order, do updating across all cases

Earlier version	[1]	1	2	3	NA	5	6	NA
later version	[1]	NA	NA	3	2	5	4	NA
Comparison		NA	NA	T	NA	T	F	NA
updated		1	2	3	2	5	4	NA

If F: override
If(earlier is NA) : override
else remain