

Supplemental Tables and Figures for The impact of data quality and source data verification on epidemiologic inference: a practical application using HIV observational data

Supplemental Table 1: Complete list of study variables with descriptions.

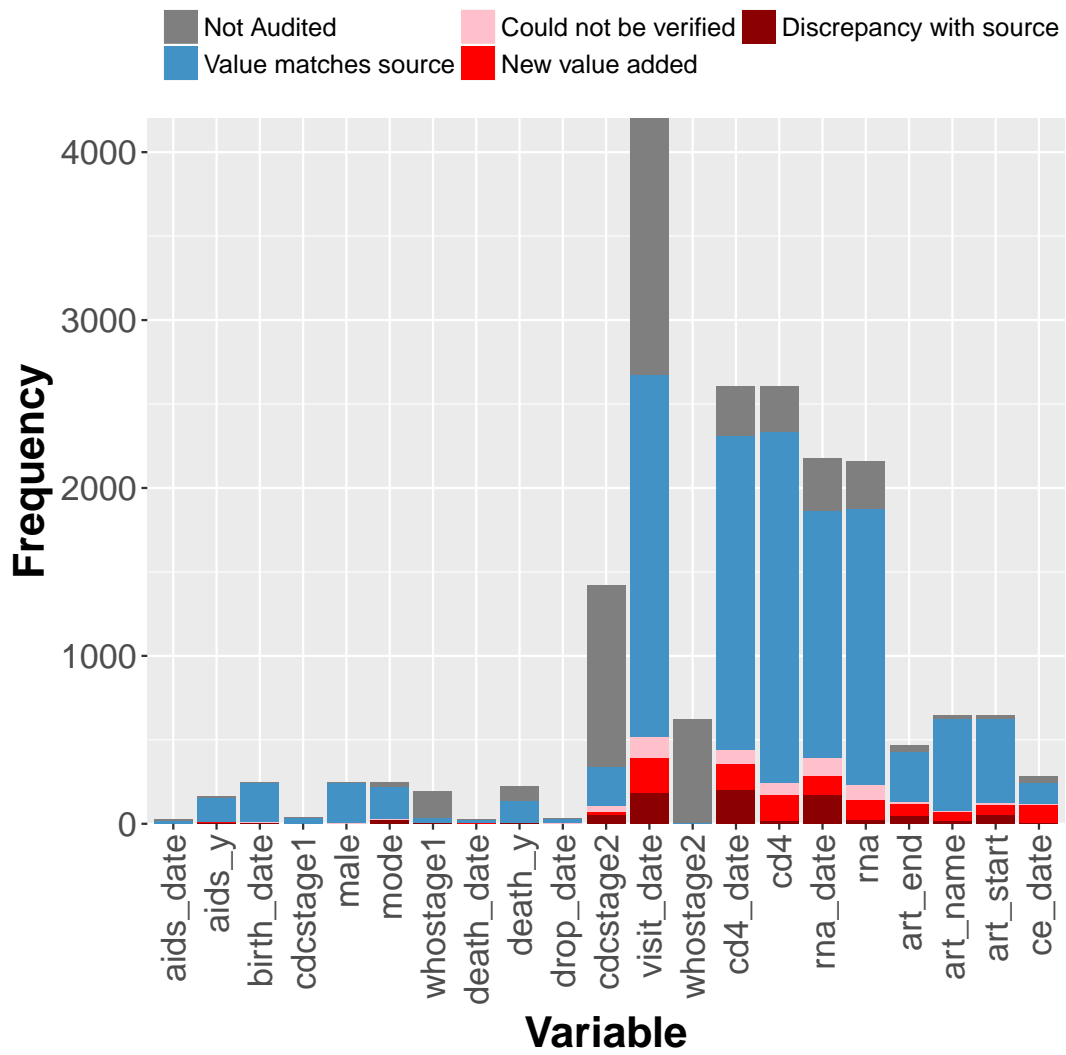
Variable	Description	Source	Dataset			Variable Audited in 2008-09
			Pre-Audit	Audited	Post-Audit	
<i>birth_date</i>	Birth date	tblBasic	Yes	Yes	Yes	Yes
<i>male</i>	Gender at birth (0-Female, 1-Male, 9-Unknown)	tblBasic	Yes	Yes	Yes	Yes
<i>mode</i>	Probable Mode of Infection	tblBasic	Yes	Yes	Yes	Yes
<i>aids_y</i>	AIDS dx before 1st visit	tblBasic	Yes	Yes	Yes	No
<i>aids_date</i>	AIDS dx date (if <i>aids_y</i> =1)	tblBasic	Yes	Yes	Yes	No
<i>cdstage</i> ^a	CDC Stage at enrollment	tblBasic	Yes	Yes	No	Yes
<i>whostage</i> ^a	WHO Stage at enrollment	tblBasic	Yes	Yes	No	Yes
<i>aids_cl_y</i> ^a	Clinical AIDS dx before 1st visit	tblBasic	No	No	Yes	No
<i>aids_cl_date</i> ^a	Clinical AIDS dx date (if <i>aids_y</i> =1)	tblBasic	No	No	Yes	No
<i>death_y</i>	Did patient pass away?	tblFollow	Yes	Yes	Yes	Yes
<i>death_date</i>	Death date (if <i>death_y</i> =1)	tblFollow	Yes	Yes	Yes	Yes
<i>drop_date</i>	Date patient dropped from cohort	tblFollow	Yes	Yes	Yes	No
<i>visit_date</i>	Clinical encounter date	tblVisit	Yes	Yes	Yes	Yes
<i>cdstage</i>	CDC Stage	tblVisit	Yes	Yes	Yes	No
<i>whostage</i>	WHO Stage	tblVisit	Yes	Yes	Yes	No
<i>cd4_date</i>	Date of CD4 lab test	tblLab.CD4	Yes	Yes	Yes	Yes
<i>cd4</i>	CD4 count value	tblLab.CD4	Yes	Yes	Yes	Yes
<i>rna_date</i>	Date of RNA lab test	tblLab.RNA	Yes	Yes	Yes	Yes
<i>rna</i>	RNA value	tblLab.RNA	Yes	Yes	Yes	Yes
<i>art_start</i>	Date of ART drug start	tblART	Yes	Yes	Yes	Yes
<i>art_end</i>	Date of ART drug end	tblART	Yes	Yes	Yes	Yes
<i>art_name</i>	Code representing ART drug(s)	tblART	Yes	Yes	Yes	Yes
<i>ce_date</i>	Date of clinical outcome	tblCE	Yes	Yes	Yes	No
<i>base_sex</i>	Gender at birth (0-Female, 1-Male, 9-Unknown)	Derived	Yes	Yes	Yes	
<i>base_mode</i>	Probable Mode of Infection	Derived	Yes	Yes	Yes	
<i>base_year</i>	Year of ART initiation	Derived	Yes	Yes	Yes	
<i>base_age</i>	Age at ART initiation	Derived	Yes	Yes	Yes	
<i>base_ADE</i>	Clinical AIDS at ART initiation	Derived	Yes	Yes	Yes	
<i>base_ARTregimen</i>	1st ART regimen	Derived	Yes	Yes	Yes	
<i>base_CD4</i>	CD4 cell count at ART initiation	Derived	Yes	Yes	Yes	
<i>base_nadirCD4</i>	Lowest CD4 cell count prior to ART initiation	Derived	Yes	Yes	Yes	
<i>base_VL</i>	Viral load at ART initiation	Derived	Yes	Yes	Yes	
<i>base_VLundetectable</i>	Undetectable VL at ART initiation	Derived	Yes	Yes	Yes	
<i>event_death</i>	Did patient pass away?	Derived	Yes	Yes	Yes	
<i>event_ADE</i>	Did patient have post-ART clinical outcome?	Derived	Yes	Yes	Yes	
<i>time_follow</i>	Time from ART initiation to death/censoring	Derived	Yes	Yes	Yes	
<i>time_ADE</i>	Time from ART initiation to ADE/censoring	Derived	Yes	Yes	Yes	

The data protocol further defining variables is available at <https://www.ccasanet.org/resources/>.

^a In 2014, the data transfer protocol pertaining to clinical AIDS status at enrollment was updated. Certain variables were deprecated in tblBasic and replaced with new variables. These variables are not included in comparisons of pre-audit and post-audit datasets.

Supplemental Table 2: Overview of audit frequency by site.

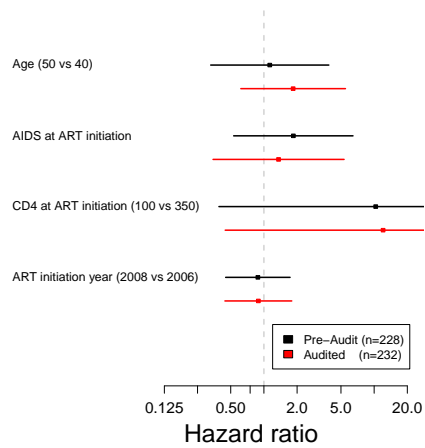
Site	Audited N(%)	Not Audited N(%)	Total
Site A	31 (89%)	4 (11%)	35
Site B	12 (57%)	9 (43%)	21
Site C	31 (52%)	29 (48%)	60
Site D	34 (97%)	1 (3%)	35
Site E	28 (80%)	7 (20%)	35
Site F	30 (100%)	0 (0%)	30
Site G	30 (100%)	0 (0%)	30
Site H	24 (69%)	11(31%)	35
Site I	30 (86%)	5 (14%)	35
Overall	250	66	316



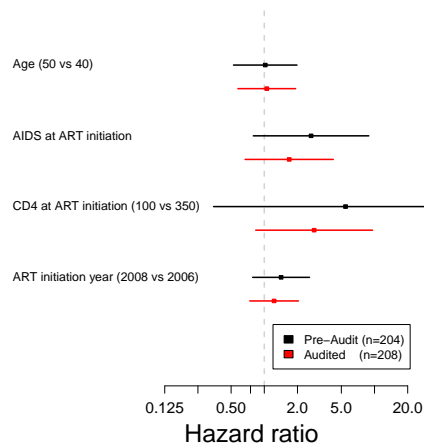
Supplemental Figure 1: Summary of reported audit findings for all audited variables.

Supplemental Table 3: Auditing results for each variable entry.

Form	Variable	Total Entries	Not audited	Total Audited	Value matches source A1	Discrepancy with source A2 + A3	New value added A4	Could not be verified A5	Error Rate $(A2+A3+A4+A5)/(A1+A2+A3+A4+A5)$
basic	<i>male</i>	250	5	245	242	2	0	1	1.2%
	<i>cdcstage1</i>	35	0	35	34	1	0	0	2.9%
	<i>birth.date</i>	250	4	246	237	7	0	2	3.7%
	<i>aids.y</i>	164	6	158	146	9	1	2	7.6%
	<i>mode</i>	250	28	222	192	23	0	7	13.5%
	<i>aids.date</i>	28	6	22	18	1	3	0	18.2%
	<i>whostage1</i>	192	157	35	27	7	0	1	22.9%
follow	<i>death.y</i>	226	88	138	131	5	0	2	5.1%
	<i>drop.date</i>	30	0	30	25	0	4	1	16.7%
	<i>death.date</i>	26	2	24	18	0	4	2	25.0%
visit	<i>whostage2</i>	623	616	7	7	0	0	0	0.0%
	<i>visit.date</i>	4203	1529	2674	2157	183	211	123	19.3%
	<i>cdcstage2</i>	1425	1084	341	232	54	16	39	32.0%
cd4	<i>cd4</i>	2607	273	2334	2090	17	155	72	10.5%
	<i>cd4.date</i>	2607	292	2315	1875	205	155	80	19.0%
rna	<i>rna.v</i>	2163	285	1878	1644	27	116	91	12.5%
	<i>rna.date</i>	2175	314	1861	1471	173	113	104	21.0%
art	<i>art.name</i>	644	18	626	546	17	56	7	12.8%
	<i>art.start</i>	645	21	624	496	56	57	15	20.5%
	<i>art.end</i>	467	36	431	298	46	73	14	30.9%
ce	<i>cc.date</i>	279	36	243	123	10	102	8	49.4%
	<i>overall</i>	19289	4800	14489	12009	843	1066	571	17.1%

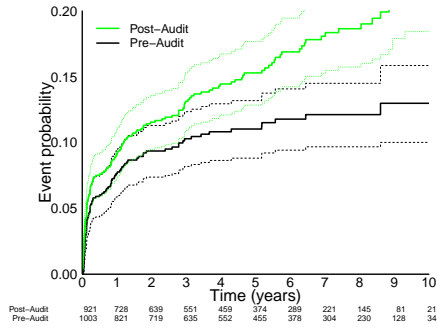


(a) Mortality

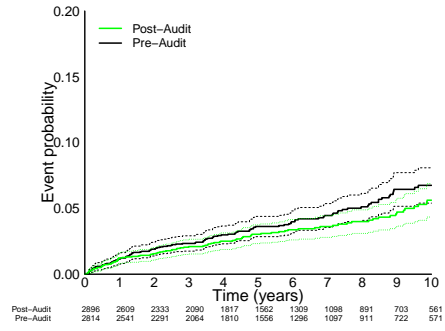


(b) AIDS-defining event

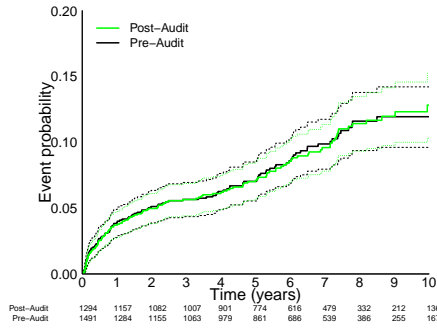
Supplemental Figure 2: Adjusted hazard ratios of mortality (a) and AIDS-defining event (b) for patients in the pre-audit and audited datasets.



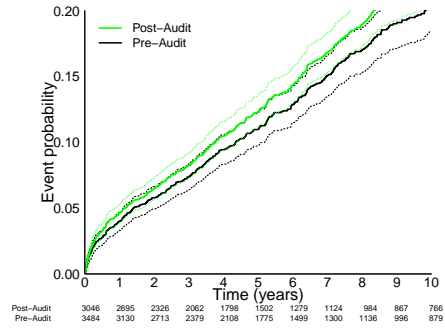
(a) Site 1



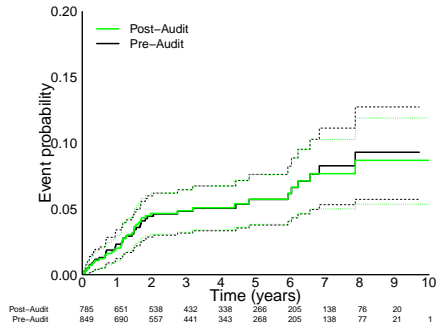
(b) Site 2



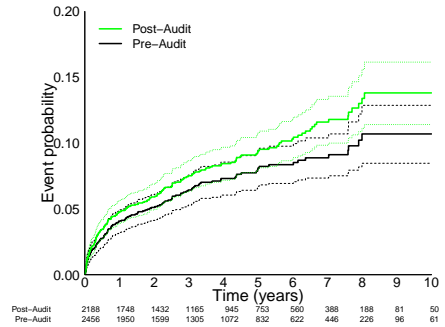
(c) Site 3



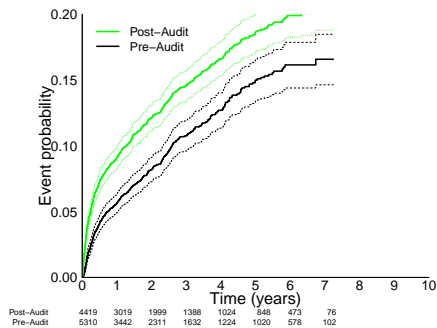
(d) Site 4



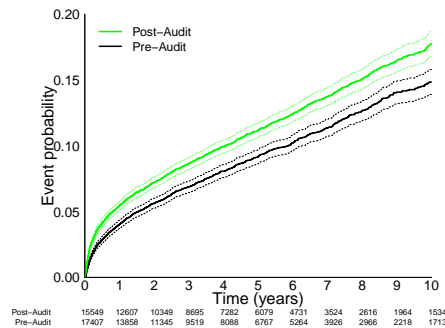
(e) Site 5



(f) Site 6

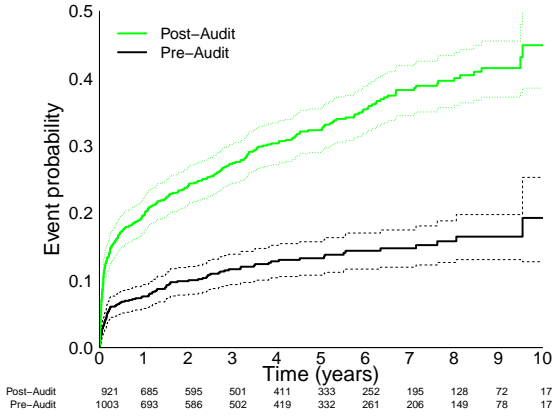


(g) Site 7

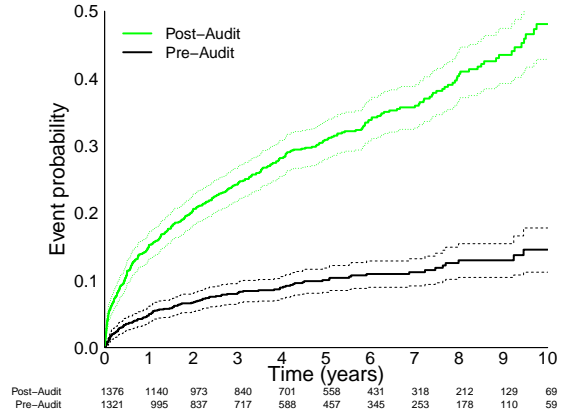


(h) Overall

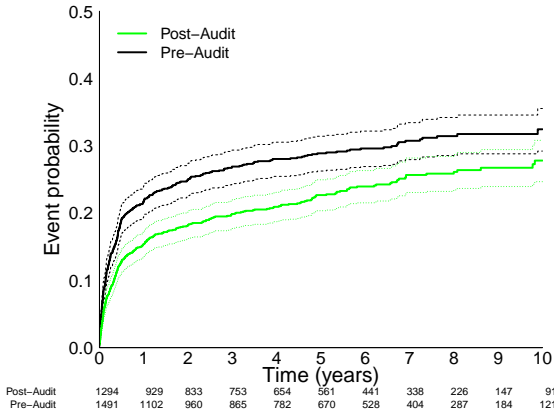
Supplemental Figure 3: Estimated cumulative incidence of death by site for patients in the pre-audit and post-audit datasets. Solid lines denote the estimated incidence and dotted lines denote the corresponding 95% confidence intervals.



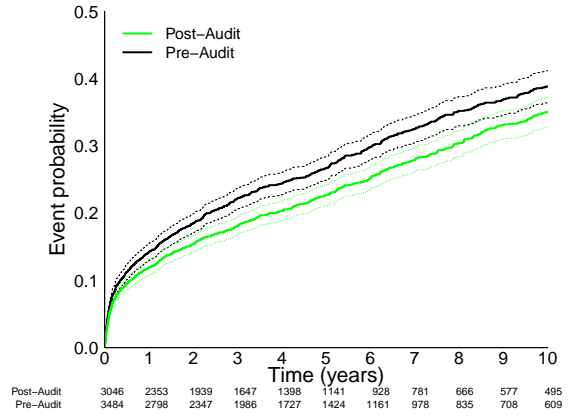
(a) Site 1



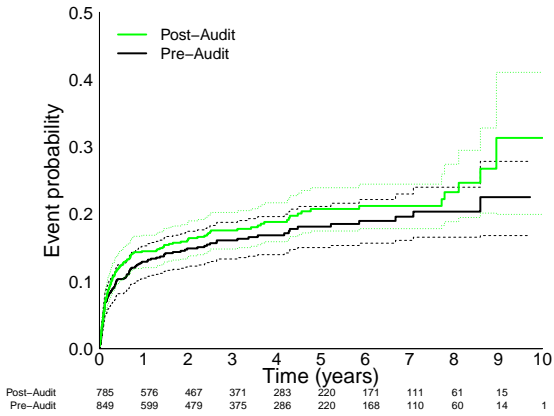
(b) Site 2



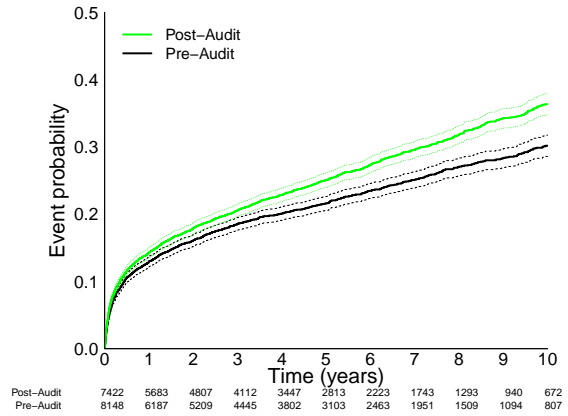
(c) Site 3



(d) Site 4

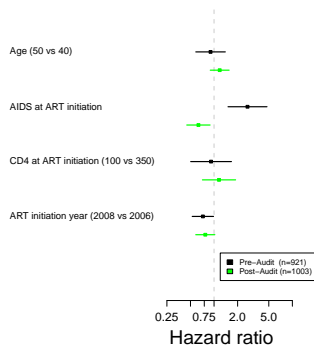


(e) Site 5

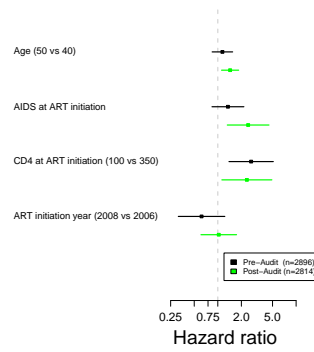


(f) Overall

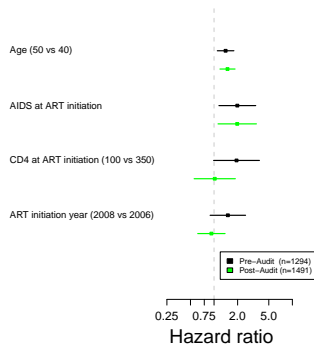
Supplemental Figure 4: Estimated cumulative incidence of an AIDS-defining event by site for patients in the pre-audit and post-audit datasets. Solid lines denote the estimated incidence and dotted lines denote the corresponding 95% confidence intervals.



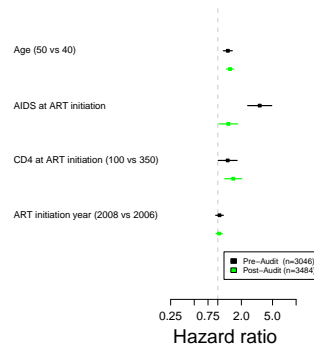
(a) Site 1



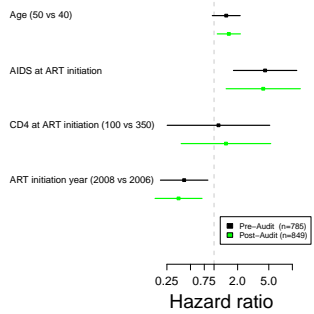
(b) Site 2



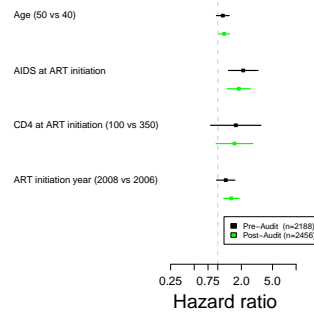
(c) Site 3



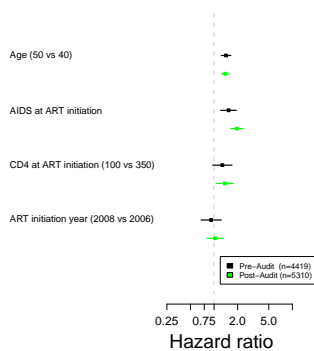
(d) Site 4



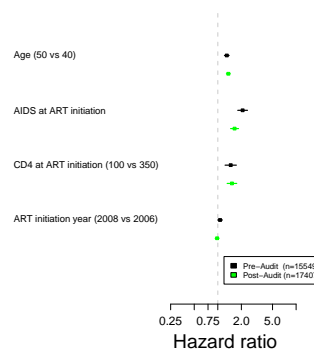
(e) Site 5



(f) Site 6

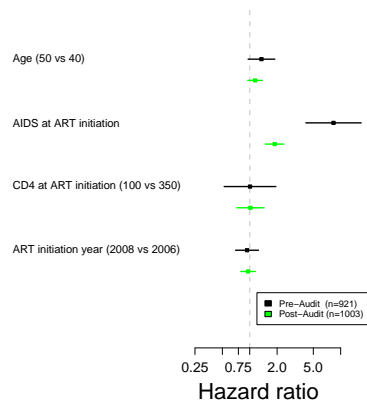


(g) Site 7

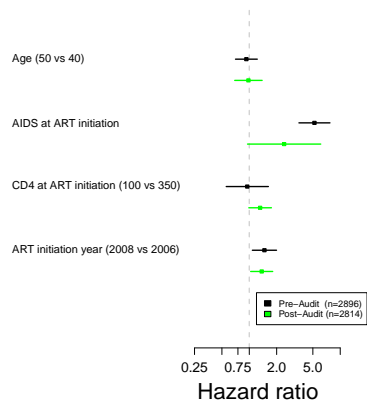


(h) Overall

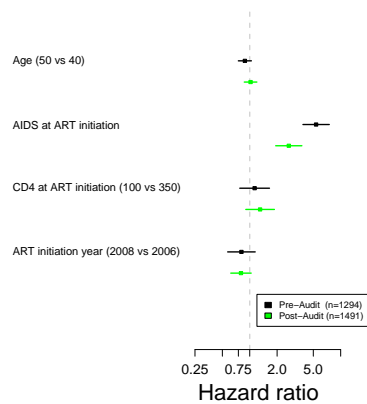
Supplemental Figure 5: Association between patient characteristics at baseline and hazard of death by site using pre-audit and post-audit data. The dots denote the estimated incidence and the lines denote the corresponding 95% confidence intervals.



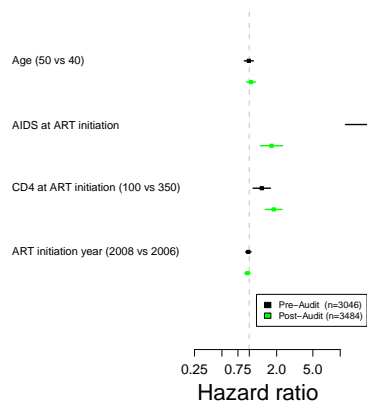
(a) Site 1



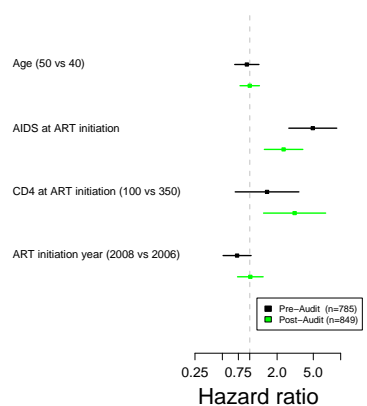
(b) Site 2



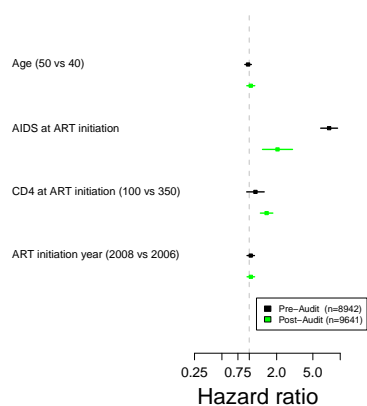
(c) Site 3



(d) Site 4



(e) Site 5



(f) Overall

Supplemental Figure 6: Association between patient characteristics at baseline and hazard of an AIDS-defining event by site using pre-audit and post-audit data. The dots denote the estimated incidence and the lines denote the corresponding 95% confidence intervals.