# HIV/AIDS Review-Verify-Extract Source Tool (HARVEST)

**Internal Tool & ETL Pipeline for HIV/AIDS Dashboard Visualizing Country Reports, PEPFAR, and UNAIDS Data**

**Ari Israel**

**Supervised by Derek Azar & Timothy Fowler**

**Population Division - Health Studies Branch**

Shape
your future
START HERE >

United States®
Census
2020

# Goals

→ Visualize HIV/AIDS pandemic in a dashboard similar to JHU COVID-19 map
→ Emphasize confirmed data made public by country governments
→ Collate public data from UNAIDS and PEPFAR in addition to country reports
→ Identify relevant sources to extract and verify statistics from annual reports
→ Monitor sources for updated reports or tables to include in database

Shape
your future
START HERE >

United States®
Census
2020

# Challenges

→ Limited universe of confirmed HIV/AIDS pandemic data
  → Most countries don't publish annual surveillance reports
  → Most data consistently available (i.e. UNAIDS) is modeled
  → Lots of broken MOH sites/dead links for developing countries

→ Available country-reported data in unstructured format i.e. PDFs

→ Identifying relevant sources and extracting statistics not easily automated

→ Extracted statistics need to be manually verified by second analyst

Shape
your future
START HERE >

United States®
Census
2020

# Solutions

→ Automated crawl/scrape of manually-identified sites with annual reports
→ Internal tool for analysts to review sources, extract and verify statistics
→ Save data collected via tool to database with ability to export master table

Shape
your future
START HERE >

United States®
Census
2020

# Technologies

→ **Anaconda** for environment/package management
→ **GitLab** for source control and wiki docs
→ **SQLite** to store prototype data
→ **Scrapy** to crawl and scrape sites
→ **Requests** to fetch data via APIs
→ **Python** for backend and Django
→ **Django** for web app framework and ORM
→ **JavaScript** for review interface frontend
→ **HTML/CSS** for Django templates and styling
→ **PDF.js** for interactive PDFs in review interface

Shape
your future
START HERE >

United States®
Census
2020

# Research and analysis

# HIV in the United Kingdom

Reports by Public Health England about testing, diagnosis and care HIV in the UK.

Published 1 November 2013
Last updated 17 January 2020 — see all updates
From: **Public Health England**

**Scrapy** used to crawl and scrape UK government site with list of reports

Results shown in HARVEST **Django** app

## Documents

HIV in the UK: towards zero HIV transmissions by 2030, 2019 report

Ref: PHE publications gateway number: GW-920
PDF, 1.91MB, 88 pages

HIV in the UK: towards zero HIV transmissions by 2030, 2019 appendix

Ref: PHE publications gateway number: GW-920
PDF, 1.03MB, 28 pages

## Sources crawled

- HIV in the United Kingdom

## Reports found

- HIV in the UK: towards zero HIV transmissions by 2030, 2019 report
- Progress towards ending the HIV epidemic in the UK: 2018 report
- Towards elimination of HIV transmission, AIDS and HIV-related deaths in the UK: 2017 report
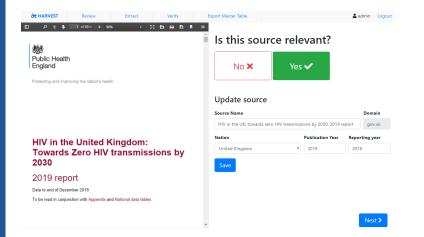- HIV in the UK: 2016 report
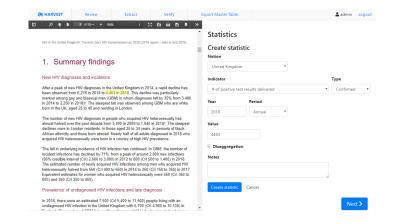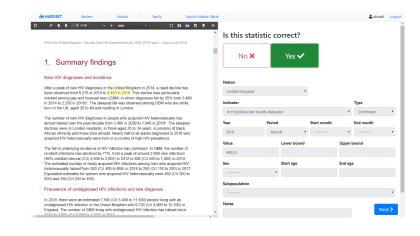
# Workflow

1. **Review** source for relevance → 2. **Extract** statistics from source → 3. **Verify** extracted statistics



Analyst 1



Analyst 1



Analyst 2

Shape
your future
START HERE >

United States®
Census
2020

# Deliverables

- ✓ **System design and app documentation**

- ✓ **HARVEST internal tool prototype**

**Next step: productionize software**
Document business logic and language-agnostic architecture for potential DSD rewrite in .NET

Shape
your future
START HERE >

United States®
Census
2020