

# Data Resource Profile: Prostate cancer data from Clinical Practice Research Datalink linked hospital records, mortality data and cancer registry standardized to the Observational Medical Outcomes Partnership common data model (CPRD-PCa-OMOP)

Eng Hooi Tan<sup>1, </sup>, Danielle Newby<sup>1</sup>, Daniel Prieto-Alhambra<sup>1,2, </sup>, Mandickel Kamtengeni<sup>1</sup>, Antonella Delmestri<sup>1,\* </sup>, OPTIMA Consortium<sup>†</sup>

<sup>1</sup>Health Data Sciences, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, United Kingdom

<sup>2</sup>Department of Medical Informatics, Erasmus University Medical Centre, Rotterdam, The Netherlands

\*Corresponding author. Nuffield Department of Orthopaedics, BOTNAR Research Centre, Rheumatology and Musculoskeletal Sciences, University of Oxford, Oxford, OX3 7LD, United Kingdom. E-mail: antonella.delmestri@ndorms.ox.ac.uk

<sup>†</sup>See Acknowledgments section for a complete list of authors who are part of the OPTIMA Consortium as well as their affiliations.

**Keywords** CPRD, National Cancer Registration and Analysis Service, NCRAS, OMOP CDM, oncology, data harmonization, data linkage

## Key Features

- CPRD-PCa-OMOP comprises men from England with incident prostate cancer (PCa) selected from the Clinical Practice Research Datalink (CPRD) GOLD and Aurum primary care databases, linked to hospital admissions, Office for National Statistics mortality data, and national cancer registry records.
- Each dataset was standardized to the Observational Medical Outcomes Partnership Common Data Model to facilitate data integration, collaborative research, external validation, and network studies. All the standardized datasets were merged to create CPRD-PCa-OMOP.
- Men aged  $\geq 18$  years, diagnosed with their first PCa in primary care between 2010 and 2022, and confirmed in the cancer registry, with  $\geq 1$  year of general practitioner registration, were included. There are 69 109 patients from 1611 English practices: 5566 from CPRD GOLD and 63 543 from CPRD Aurum.
- CPRD-PCa-OMOP includes information on demographics, diagnoses, medications, procedures, referrals, test results, hospitalizations, episodes of care, cancer diagnoses, tumour characteristics such as Tumour, Node, Metastasis staging and Gleason score, date of death, and primary cause of death.
- This comprehensive data resource will enable the study of risk factors, treatment patterns, and longitudinal outcomes in PCa. CPRD data access can be requested at <https://www.cprd.com/access-data>.

## Data resource basics

### Observational data for cancer patients in the UK

The digital collection of UK patient care data was started in the late 1980s to manage general practitioner (GP) practices and patients, and was expanded over the years to include hospital events and registry records, generating large data sources that are now used for observational research [1].

Observational data for cancer patients are complex, distributed across different care settings. This requires the integration of various pieces of information, such as tumour histology, biomarkers, and treatment regimens, as well as longitudinal

comorbidities and comedications. In non-cancer observational research, cohorts can be defined by the absence or presence of a health condition. The additional challenges of conducting observational research in cancer arise from diagnostic attributes (cancer staging, histology, grade, and biomarkers) that can determine the treatment options and prognoses [2].

### Clinical Practice Research Datalink

One of the UK data providers that allow cancer data integration through a linked data service is the Clinical Practice Research Datalink (CPRD). CPRD is a not-for-profit, cost-recovery UK government research service delivered by the Medicines and Healthcare products Regulatory Agency (MHRA) with support

Received: 20 June 2025. Accepted: 20 January 2026

© The Author(s) 2026. Published by Oxford University Press on behalf of the International Epidemiological Association. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

from the National Institute for Health and Care Research, as part of the Department of Health and Social Care [3]. CPRD hosts two primary care observational databases: CPRD GOLD [4] and CPRD Aurum [5], hereinafter called GOLD and Aurum, respectively. GOLD contains records collected by using the Vision<sup>®</sup> software from historical contributing practices in England, and from historical and current contributing practices in Wales, Scotland, and Northern Ireland. Aurum contains records collected by using the EMIS Web<sup>®</sup> software from historical and current contributing practices in England only. CPRD can provide a variety of linked patient-level data from England for both GOLD and Aurum, including the Hospital Episode Statistics (HES) Admitted Patient Care (APC), Office for National Statistics (ONS) Mortality, and National Cancer Registration and Analysis Service (NCRAS) datasets [6]. These datasets have been utilized individually or as part of a linkage in numerous peer-reviewed publications [7].

## CPRD-PCa-OMOP

CPRD-PCa-OMOP is a new data source created to characterize and assess long-term outcomes in patients with prostate cancer (Pca) in the context of the public-private partnership OPTIMA (Optimal Treatment for Patients with Solid Tumours in Europe Through Artificial Intelligence). The OPTIMA consortium aims to achieve an interoperable real-world data and evidence generation platform to improve care for patients with prostate, breast, and lung cancer. To create this data source, we used the GOLD [4, 8] and Aurum [5, 9] December 2023 releases, both linked to the latest available HES APC [10–12], NCRAS cancer registration (NCRASCR) [13–15], and ONS Mortality datasets [16–18].

To avoid possible patient duplications, we have excluded from GOLD those practices that migrated from the Vision<sup>®</sup> software to the EMIS Web<sup>®</sup> software within CPRD.

Because of the different data structures of these datasets and to facilitate the interoperability and execution of distributed network studies, we transferred each of these datasets to the Observational Medical Outcomes Partnership (OMOP) Common Data Model (CDM) version 5.4 [19] and merged the resulting databases. The OMOP CDM is a well-established standardization method, which has been adopted to facilitate interoperability and federated analytics in various frameworks such as the European Health Data and Evidence Network (EHDEN) consortium [20], the Data Analysis and Real World Interrogation Network (DARWIN) EU<sup>®</sup> initiative [21], and the OPTIMA consortium [22].

CPRD-PCa-OMOP is the first data source in OMOP CDM format that includes CPRD primary care data linked to HES APC, ONS Mortality, and NCRASCR datasets.

## Data collected

CPRD-PCa-OMOP includes 69 109 men from England aged  $\geq 18$  years diagnosed with their first Pca in GOLD or Aurum between 2010 and 2022 and confirmed in NCRASCR, with  $\geq 1$  year of registration at a GP practice [23]. In CPRD, only consenting GP practices from England are eligible for linkage and we linked primary care records to HES APC, ONS Mortality, and NCRASCR data (Table 1). Patients who opt out of having their data utilized in

clinical research were excluded from CPRD database releases. We excluded ONS Mortality- and HES-linked data that were provided with a weak linkage likelihood (1 to 5 = the most likely, the least likely), following a previously published algorithm [24]. We only accepted records with a linkage likelihood of between 1 and 2 (i.e. `hes_patient.match_rank <= 2`). To prevent data duplication, we retained only the most recently registered patient identifier when multiple identifiers were linked to the same individual via a general identifier provided in the HES data (i.e. `hes_patient.hes_gen_id`).

Table 2 reports the patient demographics of CPRD-PCa-OMOP. There were 5566 patients in GOLD and 63 543 in Aurum. The age distribution was similar between both data sources, with most patients (38.40%) first diagnosed between 70 and 79 years of age. The median number of follow-up years from date of diagnosis until the end of observation was 6.13 [interquartile range (IQR) 3.71–8.68].

In Table 3, we describe the patient characteristics of CPRD-PCa-OMOP in terms of cancer stage and grade, and the prevalence of comorbidities and comedications. Most patients had Stage 1 cancer (26.16%). A higher proportion of patients had Stage 3 and 4 cancers in Aurum as compared with GOLD. Prevalent comorbidities before Pca diagnosis included hypertension (52.30%), hyperlipidaemia (24.77%), and type 2 diabetes (14.92%). Common medications prescribed in the year before Pca diagnosis included lipid-modifying agents (45.62%), drugs for acid-related disorders (33.68%), and drugs for obstructive airway diseases (20.34%).

## Data standardization to the OMOP CDM

We standardized all the datasets described above to the OMOP CDM v. 5.4, the latest version recommended at the time of writing, to enable research collaboration through standardized analytical packages across database partners using a distributed federated network strategy. Although other CDMs exist, the

**Table 1** Key details about CPRD-PCa-OMOP.

Countries covered	England
Data origin	CPRD Aurum, CPRD GOLD, HES APC, ONS Mortality, NCRAS
Who is included?	69 109 male patients with first Pca between 2010 and 2022 from 1611 practices (1420 in CPRD Aurum, 191 in CPRD GOLD, 18.33% of all practices from England)
What is recorded?	Demographics, diagnoses, symptoms, prescriptions, referrals, test results, hospitalizations and episodes of care, cancer diagnoses and cancer treatments, date of death, and primary cause of death
Start and end dates	From January 1995 to November 2023, with a median follow-up of 6.12 years from Pca diagnosis (IQR: 3.71–8.68)

IQR, interquartile range.

**Table 2** CPRD-PCa-OMOP patient demographics.

Demographics	GOLD	Aurum	All
Number of patients	5566	63 543	69 109
<b>Age at diagnosis (years) [n (%)]</b>			
18–49	45 (0.81)	599 (0.94)	644 (0.93)
50–59	511 (9.18)	6330 (9.96)	6841(9.90)
60–69	1848 (33.20)	20 647 (32.49)	22 495 (32.55)
70–79	2082 (37.41)	24 454 (38.48)	26 536 (38.40)
80+	1080 (19.40)	11 513 (18.12)	12 593 (18.22)
<b>Ethnicity [n (%)]</b>			
Asian	12 (0.22)	253 (0.40)	265 (0.38)
Asian Indian/ Bangladeshi/Pakistani	46 (0.83)	773 (1.22)	819 (1.19)
Black	99 (1.78)	2463 (3.88)	2562 (3.71)
Chinese	8 (0.14)	89 (0.14)	97 (0.14)
White	5265 (94.59)	57 522 (90.52)	62 787 (90.85)
No matching concept	136 (2.44)	2443 (3.84)	2579 (3.73)
<b>Follow-up</b>			
Follow-up from diagnosis (years) [median (IQR)]	4.92 (2.65–7.14)	6.24 (3.88–8.83)	6.13 (3.71–8.68)
<b>Death [n (%)]</b>			
Death (all-cause)	1979 (35.56)	23 251 (36.59)	25 230 (36.51)
Death (PCa)	874 (15.70)	9290 (14.62)	10 164 (14.71)

IQR, interquartile range.

**Table 3** CPRD-PCa-OMOP patient clinical characteristics.

Characteristics	GOLD patients	Aurum patients	All patients
Number of patients	5566	63 543	69 109
<b>Cancer stage [n (%)]</b>			
1	1430 (25.69)	16 650 (26.20)	18 080 (26.16)
2	982 (17.64)	10 248 (16.13)	11 230 (16.25)
3	816 (14.66)	11 207 (17.64)	12 023 (17.40)
4	831 (14.93)	10 769 (16.95)	11 600 (16.79)
Missing	1507 (27.08)	14 669 (23.09)	16 176 (23.41)
<b>Gleason combined grade [n (%)]</b>			
Low (2–6)	1210 (21.74)	12 707 (20.00)	13 917 (20.14)
Medium (7)	2225 (39.97)	25 407 (39.98)	27 632 (39.98)
High (8–10)	1204 (21.63)	15 083 (23.74)	16 287 (23.57)
Invalid <sup>a</sup> /missing	927 (16.65)	10 346 (16.28)	11 273 (16.31)
<b>Charlson Comorbidity Index (2 years before PCa diagnosis) [n (%)]</b>			
0	4764 (85.59)	52 939 (83.31)	57 703 (83.50)
1	454 (8.16)	5702 (8.97)	6156 (8.91)
2	178 (3.20)	2649 (4.17)	2827 (4.09)
≥3	144 (2.59)	2082 (3.28)	2226 (3.22)
Missing	26 (0.47)	171 (0.27)	197 (0.29)
<b>Comorbidities any time before PCa diagnosis [n (%)]</b>			
Hypertension	2563 (46.05)	33 582 (52.85)	36 145 (52.30)
Hyperlipidaemia	1301 (23.37)	15 815 (24.89)	17 116 (24.77)
Type 2 diabetes	760 (13.65)	9554 (15.04)	10 314 (14.92)
<b>Comedications in the year before PCa diagnosis [n (%)]</b>			
Lipid-modifying agents	2477 (44.50)	29 049 (45.72)	31 526 (45.62)
Drugs for acid-re- lated disorders	1851 (33.26)	21 428 (33.72)	23 279 (33.68)
Drugs for obstructive air- way diseases	1140 (20.48)	12 919 (20.33)	14 059 (20.34)

<sup>a</sup> Invalid values were any values outside of 2–10.

OMOP CDM has been chosen by health regulatory agencies to support policy-making [21] and is widely used nationally [3, 25] and internationally [26]. The harmonization of the dictionaries used in the datasets was produced by the EHDS and Observational Health Data Sciences and Informatics (OHDSI) initiatives: standardized vocabularies are available for downloading at <https://athena.ohdsi.org>. For this study, we used the Athena vocabulary bundle v5.0 29-FEB-24 and the sources for standard mapped values are documented in [Supplementary Tables S1–S11](#).

The data standardization to the OMOP CDM for primary care, secondary care, and mortality data was performed by following previously used and established techniques ([https://github.com/oxford-pharmacoepi/etl\\_ndorms](https://github.com/oxford-pharmacoepi/etl_ndorms)). For the cancer registration data, we developed new coding techniques by following the approach recommended by the OHDSI Oncology Working Group [2]. This approach, which provided part of the supporting rationale for the development of the OMOP CDM version 5.4, recognizes that cancer observational data are more challenging than those for most other conditions because they contain an individual's sequence of cancer states and treatments that define the journeys of patients through the disease. Because CPRD-PCa-OMOP includes ONS Mortality data, which are based on death certificates and are the primary source of information for NCRAS mortality, we did not utilize the NCRAS information on death. For all the datasets, source data field names and values were stored in the corresponding source reference variables in the CDM tables to maintain data provenance. [Table 4](#) summarizes the information available in the linked datasets and the corresponding source and target vocabularies.

### Cancer diagnoses

To map NCRAS cancer diagnoses of the 'tumour' dataset, we employed the World Health Organization International Classification of Diseases for Oncology, 3rd Edition (ICD-O-3) international standard for tumour registries, which uses a concatenation of cancer's attributes: [Histology]/[behaviour]-[Topography] ([Supplementary Table S1](#)). We mapped 69 106 PCA diagnoses to the OMOP CDM CONDITION\_OCCURRENCE table by using 26 standard concept identifiers, of which 24 belonged to the ICD-O-3 vocabulary. One of the two standard concept identifiers that belonged to the Systemized Nomenclature of Medicine (SNOMED) vocabulary (4161028 = '8140/3-C61.9' = 'Adenocarcinoma of prostate') covered 92.1% of the diagnoses. Six Aurum patients had diagnoses not present in the OHDSI vocabularies: the latter were stored in the OMOP CDM OBSERVATION table.

### Cancer modifiers

Whenever possible, we mapped cancer modifiers of the 'tumour' dataset (i.e. stage, grade, size, number of local and regional lymph nodes tested and those classified as metastatic, and Charlson comorbidities indexes over 2 and 6 years) into the OMOP CDM MEASUREMENT table by using the OHDSI 'Cancer Modifier' and SNOMED standardized vocabularies, for a total of 510 359 records. When cancer modifiers were not adequately represented in the OHDSI vocabularies, we stored them in the OMOP CDM OBSERVATION table, following the OHDSI best practice recommendations [27].

### Cancer staging

Using the 'Cancer Modifier' vocabulary, we mapped several cancer stage attributes reported in the NCRAS dataset and classified them by using the American Joint Committee on Cancer (AJCC)/Union for International Cancer Control (UICC) Tumour, Node, Metastasis (TNM) system [28]. AJCC/TNM is a globally recognized standard that categorizes malignant tumours based on the primary tumour (T), the regional lymph node involvement (N), and distant metastasis (M). NCRAS derived stage attributes by combining related variables: 'stage\_best' ('t\_best', 'n\_best', 'm\_best'), 'stage\_img' ('t\_img', 'n\_img', 'm\_img'), and 'stage\_path' ('t\_path', 'n\_path', 'm\_path'). We used 43 standard concept identifiers to map stages 1[A–E] to 4[A–E] associated with different versions of the AJCC/UICC TNM staging system for different patients and diagnoses (i.e. 6th, 7th, 8th, and no version).

### Cancer grading

We mapped to standard concept identifiers several cancer grade attributes: the 'Grade of tumour' (G1 to G4 = more abnormal cells, greater likelihood of aggressive growth and spread) ([Supplementary Table S2](#)) and the scoring system for malignant neoplasm of the prostate 'Gleason Primary/Secondary/Tertiary Pattern Grade' (1 to 5 = the more aggressive, the least differentiated) by using the 'Cancer Modifier' vocabulary ([Supplementary Tables S3–S5](#)), and the 'Gleason Combined Pattern Grade' ([Supplementary Table S6](#)) by using the SNOMED vocabulary.

There were 28 088 'Grade of tumour' instances, of which 51.08% were classified as G3, 32.94% as G2, 12.49% as G4, and 3.50% as G1.

The 57 852 'Gleason Combined' instances were associated with 21 distinct values, of which we accepted values from 2 to 10. The most frequent value was 7 (39.98%), followed by 6 (19.94%), 9 (13.36%), 8 (9.07%), and 10 (1.13%).

### Cancer size

There were 440 instances of 'Size of tumour' reported as the diameter of a tumour in millimetres (if more than one tumour was present, the size of the largest tumour was reported). These instances were associated with 66 distinct values, of which the most frequent 5 (20, 15, 25, 10, and 12 mm) represented 33.18% of the total.

### Cancer lymph nodes

There were 4482 instances of 'Number of local and regional metastatic lymph nodes'. Of these, 17.0% reported no positive lymph nodes (value = 0) and 50.4% reported values from 1 (8.41%) to 9 (3.10%). The remaining 32.62% covered values from 11 (3.10%) to 89 (0.02%).

There were 4375 instances of 'Number of local and regional lymph nodes examined' in excised specimens, with 89.74% reporting a value of 0.

### Charlson comorbidities indexes

There were 68 912 instances of both Charlson comorbidity scores on diagnosis date: one covering 2 years, looking back from 27 to 3 months before diagnosis; and the other covering 6 years, looking back from 78 to 6 months before diagnosis. In both cases, the most frequent score was 0 (83.50% and 74.83%,

**Table 4** Key information and vocabularies of CPRD-PCa-OMOP.

Dataset	Key information	Source vocabulary	OMOP table	Target vocabulary
CPRD	Diagnoses	Read SNOMED	condition_occurrence	SNOMED, OMOP Extension
	Procedures	Read SNOMED	procedure_occurrence	SNOMED, CPT4, OMOP Extension
	Medications	dm+d gemscript	drug_exposure	RxNorm, RxNorm Extension CVX
	Medical devices	dm+d gemscript	device_exposure	dm+d, SNOMED
	Measurements	Read SNOMED	measurement	SNOMED, LOINC, OMOP Extension
	Anything else	Read SNOMED	observation	SNOMED + many
HES APC	Diagnoses	ICD-10	condition_occurrence	SNOMED
	Procedures	OPCS-4	procedure_occurrence	SNOMED
ONS	Primary cause of death	ICD-10	death	SNOMED
NCRASCR	Cancer diagnoses	ICD-O-3 + many	condition_occurrence	SNOMED, ICD-O-3
	Cancer modifiers	Bespoken	measurement	Cancer Modifier, SNOMED
	Cancer procedures	OPCS-4	procedure_occurrence	SNOMED, OPCS-4
	Cancer medications	Bespoken RxNorm Extension	drug_exposure	RxNorm, RxNorm Extension

CPT4, Current Procedural Terminology, Fourth Edition; CVX, vaccine administered code; dm+d, NHS dictionary of medicines and devices; ICD-10, International Statistical Classification of Diseases and Related Health Problems, 10th Revision; ICD-O-3, International Classification of Diseases for Oncology, 3rd Edition; LOINC, Logical Observation Identifiers Names and Codes; NCRASCR, National Cancer Registration and Analysis Service cancer registration data; OPCS-4, Classification of Interventions and Procedures, version 4; SNOMED, Systemized Nomenclature of Medicine.

respectively), followed by 1 (8.91%, 13.49%), 2 (4.09%, 6.30%), 3 (1.80%, 2.80%), 4 (0.86%, 1.28%), and 5 (0.36%, 0.59%), with higher scores having minimum representation.

### Cancer procedures

Cancer procedures from the NCRAS 'Treatment' dataset were mapped to the OMOP CDM PROCEDURE\_OCCURRENCE table by using the OHDSI SNOMED and OPCS-4 standardized vocabularies (Supplementary Tables S7 and S8), for a total of 159 673 records. The five most frequent procedures accounted for 81.64% of the records, with 'Imaging' covering 40.07%, 'Hormone therapy' covering 20.27%, 'External beam radiation therapy procedure' covering 9.66%, 'Transrectal needle biopsy of prostate' covering 6.14%, and 'Radical prostatectomy' covering 5.50%.

### Cancer treatments

Cancer treatments from the NCRAS 'Treatment' dataset were mapped to the OMOP CDM DRUG\_EXPOSURE table by using OHDSI 'RxNorm' and 'RxNorm Extension' standardized vocabularies (Supplementary Table S9). A total of 2913 records were mapped using 39 standard concept identifiers and the five most frequent concepts covered 85.65% of the records. Bicalutamide

accounted for 38.76%, goserelin for 26.16%, cyproterone for 14.07%, gonadorelin for 3.50%, and docetaxel for 3.16%.

### Cancer episodes

All cancer diagnoses, modifiers, and treatments (Supplementary Table S10) were additionally mapped to the OMOP CDM EPISODE and EPISODE\_EVENT tables to identify cancer patients' trajectories through the disease (Supplementary Table S11). A total of 243 095 episodes were created associated with the domains of Procedure (42.61%), Regimen (28.96%), and Condition (28.43%). The number of episode events was 1 548 516, linking the events to the OBSERVATION (46.40%), MEASUREMENT (38.64%), PROCEDURE\_OCCURRENCE (10.31%), CONDITION\_OCCURRENCE (4.46%), and DRUG\_EXPOSURE (0.19%) tables.

### Cancer deaths

All primary causes of death were mapped to the OMOP CDM DEATH table, together with the date of death. A total of 25 230 deaths were reported, with 40.29% (39.96% in Aurum, 44.16% in GOLD) associated with 'Primary malignant neoplasm of prostate' (concept\_id = 200962).

## Data-quality assessment

We assessed the quality of the data mapped to the OMOP CDM by using Kahn's framework [29], implemented by the OHDSI software DataQualityDashboard (DQD, version: 2.6.3), available at <https://github.com/OHDSI/DataQualityDashboard>. DQD performs validation and verification checks according to conformance, completeness, and plausibility at the concept, field, and table levels by using check-specific thresholds. A total of 4684 data-quality checks were performed on the data, of which 98% passed overall, with both conformance and completeness satisfying 99% of the quality assessment and plausibility 96%. The data-quality checks that failed were investigated and none of them was a cause for concern for our study (e.g. `unit_concept_ids` used in the UK were not recognized by DQD and a few checks contained software mistakes, all of which we reported to the DQD developers).

## Data resource use

The linked datasets in their source format have been used in other studies [30] to investigate the quality and concordance of cancer diagnoses. CPRD-PCa-OMOP is a novel dataset mapped to the CDM to allow participation in multinational network cohort studies. The data resource was created to contribute data to the OPTIMA consortium and it will be used to answer questions prioritized by healthcare professionals/patient stakeholder groups, related to the characterization of patients with PCa, risk stratification, treatment patterns, and comparative effectiveness/safety analyses [31]. We are further investigating the characteristics and outcomes of patients who received prostatectomy versus radiotherapy in a separate study [32].

## Strengths and weaknesses

The key strength of CPRD-PCa-OMOP is its data sources, which are comprehensive and representative of the population from England, with longitudinal follow-up information and linkage to hospital admissions, cancer registry data, and mortality data, the latter being particularly important for the accuracy of cancer survival prediction. Having linked data sources allows researchers to triangulate between data points and define the variables according to the study objectives. For example, the first date of cancer diagnosis varied between the different data sources. Whitfield *et al.* [33] found that agreement in diagnosis dates was variable across different cancer sites and showed higher agreement between HES and NCRAS. For the description of this data source, we used the date recorded in NCRAS as the reference standard. A small proportion of patients, however, had earlier dates in primary care and HES records. As noted by Whitfield *et al.* [33], earlier dates in primary care may reflect the coding of suspected cancers before confirmation. Similarly, earlier cancer diagnoses for a minority of patients in HES records may have resulted from delays in NCRAS recording. NCRAS records rely on pathological confirmation, which is a time-consuming process, whereas HES records capture the first clinically relevant diagnosis, such as one based on imaging. Another major strength of CPRD-PCa-OMOP is that it is OMOP-standardized, which

facilitates data integration, collaborative research, external validation, and network studies. Mapping of cancer records has also been performed by Genomics England [34], but we have additionally mapped episodes to study patient trajectories and followed the OMOP CDM best practice [19].

The main weakness of CPRD-PCa-OMOP is the lack of information about chemotherapy and radiotherapy regimens, which we requested after obtaining ethical approval, but never received due to a delay in the release of the datasets held by a third party. This situation has limited the capacity of the study to investigate disease-free survival, relapse or recurrence, and cancer progression. However, we remain committed to standardizing the missing regimen data and releasing the corresponding GitHub code and documentation to the research community as soon as is feasible.

Some information regarding cancer stage (23.41%) and Gleason grade (16.31%), which are important prognostic factors in cancer research, was missing. However, Strongman *et al.* [30] showed that the completeness of cancer stage and grade has improved over time, especially since 2012.

Finally, data standardization could sometimes hinder interoperability between studies when different standard concepts are used to represent the same variables. To reduce this risk, we have provided detailed documentation of the standardization process along with the corresponding GitHub code. In addition, while concept harmonization is based on the evolving Athena vocabularies, their changes mainly include additions and error corrections, ensuring that the final outputs remain highly compatible.

## Data resource access

Source data documentation are available at <https://www.cprd.com/primary-care-data-public-health-research> and <https://www.cprd.com/cprd-linked-data>. The Extract, Transform, and Load (ETL) process is documented at [https://github.com/oxford-pharmacoepi/etl\\_ndorms](https://github.com/oxford-pharmacoepi/etl_ndorms) and [https://oxford-pharmacoepi.github.io/etl\\_ndorms/docs/NCRASCR](https://oxford-pharmacoepi.github.io/etl_ndorms/docs/NCRASCR). Researchers can apply for CPRD data access subject to licensing fees and protocol approval (enquiries: [enquiries@cpd.com](mailto:enquiries@cpd.com)). For queries about the mapped dataset, please refer to the 'Supplementary Material' file, which provides detailed information on the standardized variables and on the harmonization choices. The data used in this study were provided under the University of Oxford multi-study CPRD licence and cannot be shared. However, the Python 3 code and the ETL documentation are freely available on GitHub and could be used to recreate or adapt the mapped database or as a template for other cancer data sources.

## Ethics Approval

CPRD has ethics approval from the UK Health Research Authority to support research using anonymized patient data and must complete an annual NHS Data Security and Protection Toolkit assessment to demonstrate that it meets the required standard for holding data securely. CPRD data are accessible only via a CPRD client approval, a CPRD licence, and a protocol ethical approval obtained via the CPRD's Research Data

Governance process. CPRD data cannot be shared with collaborators not named in the protocol and they need to be erased when the protocol research question has been answered. Access to CPRD data and to the linkages used to create CPRD-PCa-OMOP was granted by protocol number 22\_001867.

## Acknowledgements

We thank Teen, Wai Yi Man, from the University of Oxford for her contribution to the data-conversion process into the OMOP CDM. This study is based in part on data from the CPRD obtained under license from the UK MHRA. The data are provided by patients and collected by the NHS as part of their care and support. The interpretation and conclusions contained in this study are those of the author/s alone.

**The OPTIMA consortium (with associations):** James N'Dow<sup>1</sup>, Emma Jane Smith<sup>1</sup>, Ellen Moyse<sup>1</sup>, Carla Bezuidenhout<sup>1</sup>, Torsten Gerriet Blum<sup>1</sup>, Peter-Paul Willemse<sup>1</sup>, Philip Cornford<sup>1</sup>, Maurice Schlie<sup>1</sup>, Alberto Briganti<sup>7</sup>, Susan Evans Axelsson<sup>1</sup>, Vasileios Sakalis<sup>1</sup>, Anders Bjartell<sup>2</sup>, Monique Roobol<sup>3</sup>, Katharina Beyer<sup>3</sup>, Lionne Venderbos<sup>3</sup>, Sebastiaan Remmers<sup>3</sup>, Giorgio Gandaglia<sup>7</sup>, Raoul Boomsma<sup>3</sup>, Carolina Testa<sup>3</sup>, Charles Auffray<sup>4</sup>, Jie Shen<sup>36</sup>, Nareen Katta<sup>36</sup>, Steven MacLennan<sup>5</sup>, Sara MacLennan<sup>5</sup>, Valerie Speirs<sup>5</sup>, Yilin Xu<sup>36</sup>, Imran Omar<sup>5</sup>, Demi McDonald<sup>5</sup>, Kelly Gray<sup>5</sup>, Lesley Anderson<sup>5</sup>, Charlotte Murray<sup>5</sup>, Dianne Brown<sup>5</sup>, Abdelali (Ali) Majdi<sup>33</sup>, Carl Steinbeisser<sup>33</sup>, Christoph Kowalski<sup>22</sup>, Mart Kals<sup>6</sup>, Nora Tabea Sibert<sup>22</sup>, Juan Gómez Rivas<sup>1</sup>, Martina Faticoni<sup>7</sup>, Greta Matteuzzi<sup>7</sup>, Stephane Lejeune<sup>17</sup>, Louise Jones<sup>8</sup>, Daan Nieboer<sup>3</sup>, Javier Téllez<sup>26</sup>, Guido Juckeland<sup>9</sup>, Sandra Garrido<sup>26</sup>, Daniel Kotik<sup>9</sup>, Tobias Vu<sup>9</sup>, Artur Yakimovich<sup>9</sup>, Torsten Bauer<sup>10</sup>, Jens Kollmeier<sup>10</sup>, Roberto Galán<sup>26</sup>, Ruben Villoria<sup>26</sup>, Tobias Sjöblom<sup>11</sup>, Chatarina Larsson<sup>11</sup>, Ivaylo Stoimenov<sup>11</sup>, Daniel Prieto-Alhambra<sup>12</sup>, Sara Khalid<sup>12</sup>, Enric Bousoño Borrull<sup>26</sup>, Laura Tur Giménez<sup>26</sup>, Mahkameh Mafi<sup>12</sup>, Soraly Hernandez<sup>26</sup>, Alvaro Morandeira Galban<sup>26</sup>, Nathaly Thiel<sup>13</sup>, Nikolaus Forgó<sup>13</sup>, Antoni Napieralski<sup>13</sup>, Martina Wimmer<sup>13</sup>, Katharina Haimbuchner<sup>13</sup>, Saskia Kaltenbrunner<sup>13</sup>, Syed Shah<sup>3</sup>, Kseniia Guliaeva<sup>13</sup>, Torsten Blum<sup>10</sup>, Rebecca Graebig-Rancourt<sup>10</sup>, Michael Bussmann<sup>9</sup>, Kevin Joubel<sup>18</sup>, Talita Duarte-Salles<sup>20</sup>, Anna Palomar<sup>20</sup>, Nadia Harbeck<sup>15</sup>, Julian Koch<sup>15</sup>, Irene Lopez<sup>20</sup>, Julia Kasprzak<sup>15</sup>, Giuseppe Curigliano<sup>14</sup>, Sindhu Naidu<sup>16</sup>, Beatrice Taurelli Salimbeni<sup>14</sup>, Zoe Soo<sup>6</sup>, Aryn Bhamani<sup>16</sup>, Emanuela Ferraro<sup>14</sup>, Mario Campone<sup>18</sup>, Jean-Sebastien Frenel<sup>18</sup>, Carmen Criscitiello<sup>14</sup>, François Bocquet<sup>18</sup>, Marion Laloue<sup>18</sup>, Bérengère Denoël<sup>18</sup>, Marie Moisy<sup>18</sup>, Julien Berry<sup>18</sup>, Alexandre Auffray<sup>18</sup>, Frank Flammanc<sup>18</sup>, Hyacinthe Boni<sup>18</sup>, Clémence Gaborieau<sup>18</sup>, Josephine Bertin<sup>18</sup>, Philippe Lambin<sup>19</sup>, Anshu Ankolekar<sup>19</sup>, Chiara Corti<sup>14</sup>, Alba Tor Roca<sup>20</sup>, Elena Dal Zotto<sup>14</sup>, Andreas Kremer<sup>27</sup>, Valérie Vaccaro<sup>21</sup>, Thomy Tonia<sup>21</sup>, Céline Genton<sup>21</sup>, Maria Quaranta<sup>27</sup>, Loic Marx<sup>27</sup>, Pippa Powell<sup>23</sup>, Clare Williams<sup>23</sup>, Manou Kooy<sup>25</sup>, Nils Christian<sup>27</sup>, Christian Seeling<sup>15</sup>, Edward Burn<sup>12</sup>, Antonella Delmestri<sup>12</sup>, Inmaculada Perea Fernández<sup>26</sup>, Danielle Newby<sup>12</sup>, Cheryl Tan<sup>12</sup>, Marcel Hartig<sup>32</sup>, Juan Miguel Auñón García<sup>26</sup>, José Carlos Barrios González<sup>26</sup>, Lynn McRoy<sup>32</sup>, Chad Barnett<sup>32</sup>, Humma Khan<sup>32</sup>, Claude Chelala<sup>8</sup>, Maryam Abdollahyan<sup>8</sup>, Christian Bauer<sup>27</sup>, Mariana Pina<sup>27</sup>, Romain Tching<sup>27</sup>, Philipp Matthies<sup>28</sup>, Corinna

Zur Bonsen-Thomas<sup>28</sup>, Larissa Tschetsch<sup>28</sup>, Francisco Pinto<sup>28</sup>, Heinz Vargas<sup>28</sup>, Matthieu Blotière<sup>29</sup>, Louise Dufлот<sup>29</sup>, David Vallas<sup>29</sup>, Patrizia Torremante<sup>30</sup>, Verena Von Scharfenberg<sup>30</sup>, Nuno Azevedo<sup>31</sup>, Alla Kolyban<sup>31</sup>, Javier Cuadrado Corz<sup>8</sup>, Waltraud Kantz<sup>32</sup>, Frederic Kube<sup>32</sup>, Amanda Matthews<sup>32</sup>, Bhakti Arondekar<sup>32</sup>, Bruno Gori<sup>32</sup>, Hagen Krüger<sup>32</sup>, Julia Ilinares<sup>32</sup>, Keith Wilner<sup>32</sup>, Lucile Serfass<sup>32</sup>, Joao Mouta<sup>35</sup>, Robert Miller<sup>32</sup>, Lara Chayab<sup>35</sup>, Milan Tsompanoglou<sup>32</sup>, Karen Godbold<sup>32</sup>, Stefan Langhammer<sup>32</sup>, Anne Adams<sup>32</sup>, Sebastian Boie<sup>32</sup>, Florian Reis<sup>32</sup>, Eugenia Schander<sup>32</sup>, Andres Metspalu<sup>6</sup>, Hessa Bakhtiar<sup>32</sup>, Stefan Krautschneider<sup>32</sup>, Neda Rajaeian<sup>32</sup>, John-Edward Butler-Ransohoff<sup>33</sup>, Santiago Villalba<sup>33</sup>, Adrian Wolny<sup>33</sup>, Lisa Schneider<sup>33</sup>, Jaak Vilo<sup>6</sup>, Raivo Kolde<sup>6</sup>, Adrian Rousset<sup>34</sup>, Ivo Cleuren<sup>34</sup>, Sandra Eketorp Sylvan<sup>34</sup>, Ellie Paintin<sup>34</sup>, Monika Pokrzepa<sup>34</sup>, Carolin Lorber<sup>35</sup>, Stefanie Morris<sup>35</sup>, Sulev Reisberg<sup>6</sup>, Marek Oja<sup>6</sup>, Camille Andre<sup>35</sup>, Tobias Schulte in den Baeumen<sup>35</sup>, Jason Hannon<sup>35</sup>, Kartick Sukumaran<sup>35</sup>, Telve Objartel<sup>6</sup>, Neal Navani<sup>16</sup>, Sam Janes<sup>16</sup>, Sean C Turner<sup>36</sup>, David Dellamonica<sup>37</sup>, Heather Moses<sup>37</sup>, Yiduo Zhang<sup>37</sup>, Christophe Dufour<sup>37</sup>, Marcus Simon<sup>37</sup>, Hassan Naqvi<sup>37</sup>, Jens Ceder<sup>37</sup>, Olga Alekseeva<sup>37</sup>, Burkhard Mueller<sup>38</sup>, Tobias Flosdorf<sup>38</sup>, Daniel Kraenz<sup>38</sup>, Anastasia Goette<sup>38</sup>, Gustaf Hedström<sup>39</sup>, Peter Hellberg<sup>39</sup>, Per-Henrik Edqvist<sup>39</sup>, Rossella Nicoletti<sup>40</sup>, Mauro Gacci<sup>40</sup>, Francesco Annunziato<sup>40</sup>, Bertrand De Meulder<sup>41</sup>, Nadezhda Knauer<sup>41</sup>, Djeanne Onthoni<sup>41</sup>, Maxence Renaud<sup>41</sup>, Miguel Taraud<sup>41</sup>

<sup>1</sup>Stichting European Urological Foundation, Arnhem, Netherlands

<sup>2</sup>Lunds Universitet, Lund, Sweden

<sup>3</sup>Erasmus Universitair Medisch Centrum Rotterdam, Rotterdam, Netherlands

<sup>4</sup>Association EISBM, Vourles, France

<sup>5</sup>The University Court of the University of Aberdeen, Aberdeen, United Kingdom

<sup>6</sup>Tartu Ulikool, Tartu, Estonia

<sup>7</sup>Universita Vita-Salute San Raffaele, Milan, Italy

<sup>8</sup>Queen Mary University of London, London, United Kingdom

<sup>9</sup>Helmholtz-Zentrum Dresden-Rossendorf EV, Dresden, Germany

<sup>10</sup>Helios Klinikum Emil von Behring GMBH, Berlin, Germany

<sup>11</sup>Uppsala Universitet, Uppsala, Sweden

<sup>12</sup>The Chancellor, Masters and Scholars of the University of Oxford, Oxford, United Kingdom

<sup>13</sup>Universitat Wien, Vienna, Austria

<sup>14</sup>Istituto Europeo Di Oncologia SRL, Milan, Italy

<sup>15</sup>Ludwig-Maximiliansuniversitaet Muenchen, Munich, Germany

<sup>16</sup>University College London, London, United Kingdom

<sup>17</sup>European Organisation for Research and Treatment of Cancer AISBL, Brussels, Belgium

<sup>18</sup>Institut De Cancerologie De L'ouest, Angers, France

<sup>18</sup>Universiteit Maastricht, Maastricht, Netherlands

<sup>20</sup>Fundacio Institut Universitari Pera La Recerca A L'atencio Primaria De Salut Jordi Gol I Gurina, Barcelona, Spain

<sup>21</sup>European Respiratory Society, Lausanne, Switzerland

<sup>22</sup>Deutsche Krebsgesellschaft Ev, Berlin, Germany

<sup>23</sup>European Lung Foundation, Sheffield, United Kingdom

<sup>24-25</sup>PNO Life Sciences and Health, Rijswijk, Netherlands

<sup>26</sup>GMV Soluciones Globales Internet SAU, Madrid, Spain

<sup>27</sup>Information Technology for Translational Medicine (ITTM) SA, Esch-sur-Alzette, Luxembourg

- <sup>28</sup>Smart Reporting GMBH, Munich, Germany  
<sup>29</sup>Owkin France, Paris, France  
<sup>30</sup>PNO Innovation GMBH, Munich, Germany  
<sup>31</sup>Ydeal.Net Software LDA, Santa Maria Da Feira, Portugal  
<sup>32</sup>Pfizer Limited, Sandwich, United Kingdom  
<sup>33</sup>Bayer Aktiengesellschaft, Leverkusen, Germany  
<sup>34</sup>Amgen, Diegem, Belgium  
<sup>35</sup>F. Hoffmann-la Roche AG, Basel, Switzerland  
<sup>36</sup>Abbvie inc, North Chicago, United States  
<sup>37</sup>Astrazeneca AB, Södertälje, Sweden  
<sup>38</sup>Mutabor Technologies GMBH, Hamburg, Germany  
<sup>39</sup>Region Uppsala, Uppsala, Sweden  
<sup>40</sup>Università di Firenze, Florence, Italy  
<sup>41</sup>BDM Consulting, Lyon, France

## Author contributions

A.D. led the OMOP data standardization. A.D. and E.H.T. ran the descriptive statistics and wrote the paper. M.K. performed the data-quality checks. All the authors critically reviewed the manuscript. A.D. will act as guarantor for the paper.

## Supplementary material

Supplementary material is available at *IJE* online.

## Conflicts of interest

The Health Data Sciences research group has received research grants from the European Medicines Agency, from the Innovative Medicines Initiative, from Amgen, Chiesi, and from UCB Biopharma; as well as consultancy or speaker fees from Astellas, Amgen, AstraZeneca, and UCB Biopharma. E.H.T. received consultancy fees from Janssen Pharmaceutica NV, outside the submitted work. All other authors declare no additional conflicts of interest.

## Funding

OPTIMA is funded through the IMI2 Joint Undertaking and is listed under grant agreement No. 101034347. IMI2 receives support from the European Union's Horizon 2020 research and innovation programme and the European Federation of Pharmaceutical Industries and Associations (EFPIA). IMI supports collaborative research projects and builds networks of industrial and academic experts to boost pharmaceutical innovation in Europe. The views communicated within are those of OPTIMA. Neither the IMI nor the European Union, EFPIA, nor any Associated Partners are responsible for any use that may be made of the information contained herein. The interpretation and conclusions contained in this study are those of the author/s alone.

## Data availability

The data underlying this article were provided by the CPRD under the multi-study licence held by the University of Oxford. These data can be requested at <https://www.cprd.com/access-data> by CPRD licence holders.

## Use of Artificial Intelligence (AI) Tools

AI was not used in collecting and/or analysing the data, producing images or graphical elements, or in writing this paper.

## References

- Jick H, Jick SS, Derby LE. Validation of information recorded on general practitioner based computerised data resource in the United Kingdom. *Bmj* 1991;**302**:766–8. <https://doi.org/10.1136/bmj.302.6779.766>
- Belenkaya R, Gurley MJ, Golozar A *et al*. Extending the OMOP common data model and standardized vocabularies to support observational cancer research. *JCO Clin Cancer Inform* 2021;**5**:12–20. <https://doi.org/10.1200/cci.20.00079>
- Medicines and Healthcare products Regulatory Agency (MHRA). *Clinical Practice Research Datalink (CPRD)*. <https://cprd.com/> (31 October 2025, date last accessed).
- Sanchez-Santos MT, Axson EL, Dedman D, Delmestri A. Data Resource Profile Update: CPRD GOLD. *Int J Epidemiol* 2025;**54**:dyaf077. <https://doi.org/10.1093/ije/dyaf077>
- Wolf A, Dedman D, Campbell J *et al*. Data Resource Profile: Clinical Practice Research Datalink (CPRD) Aurum. *Int J Epidemiol* 2019;**48**:1740–g. <https://doi.org/10.1093/ije/dyz034>
- Medicines and Healthcare products Regulatory Agency (MHRA). *CPRD linked data*. <https://www.cprd.com/cprd-linked-data> (31 October 2025, date last accessed).
- Medicines and Healthcare products Regulatory Agency (MHRA). *CPRD Bibliography*. <https://www.cprd.com/bibliography> (31 October 2025, date last accessed).
- Clinical Practice Research Datalink. *CPRD GOLD December 2023 dataset*. <https://doi.org/10.48329/30pm-xq61> (31 October 2025, date last accessed).
- Clinical Practice Research Datalink. *CPRD Aurum December 2023 dataset*. <https://doi.org/10.48329/7njs-8a57> (31 October 2025, date last accessed).
- Clinical Practice Research Datalink. *CPRD GOLD HES APC January 2022*. <https://doi.org/10.48329/fagm-ez75> (31 October 2025, date last accessed).
- Clinical Practice Research Datalink. *CPRD Aurum HES APC January 2022*. <https://doi.org/10.48329/vagx-9d96> (31 October 2025, date last accessed).
- Herbert A, Wijlaars L, Zylbersztejn A, Cromwell D, Hardelid P. Data Resource Profile: Hospital Episode Statistics Admitted Patient Care (HES APC). *Int J Epidemiol* 2017;**46**:1093–i. <https://doi.org/10.1093/ije/dyx015>

13. Clinical Practice Research Datalink. *CPRD GOLD Cancer Registration August 2021*. <https://doi.org/10.48329/541y-nh70> (31 October 2025, date last accessed).
14. Clinical Practice Research Datalink. *CPRD Aurum Cancer Registration August 2021*. <https://doi.org/10.48329/5sm7-3209> (31 October 2025, date last accessed).
15. Henson KE, Elliss-Brookes L, Coupland VH *et al*. Data Resource Profile: National Cancer Registration Dataset in England. *Int J Epidemiol* 2020;**49**:16–h. <https://doi.org/10.1093/ije/dyz076>
16. Clinical Practice Research Datalink. *CPRD GOLD ONS deaths January 2022*. <https://doi.org/10.48329/WQ7V-X832> (31 October 2025, date last accessed).
17. Clinical Practice Research Datalink. *CPRD Aurum ONS deaths data January 2022*. <https://www.cprd.com/cprd-aurum-ons-deaths-data-january-2022> (31 October 2025, date last accessed).
18. Office for National Statistics (ONS). *User guide to mortality statistics*. <https://www.ons.gov.uk/peoplepopulationandcommunity/birthsdeathsandmarriages/deaths/methodologies/userguidetomortalitystatisticsjuly2017> (31 October 2025, date last accessed).
19. Observational Health Data Sciences and Informatics (OHDSI). *OMOP Common Data Model (version 5.4)*. <https://ohdsi.github.io/CommonDataModel/cdm54.html> (31 October 2025, date last accessed).
20. European Health Data and Evidence Network (EHDEN). *EHDEN Consortium*. <https://www.ehden.eu/consortium-partners/> (31 October 2025, date last accessed).
21. European Medicines Agency. *Data Analysis and Real World Interrogation Network (DARWIN EU)*. <https://www.ema.europa.eu/en/about-us/how-we-work/data-regulation-big-data-other-sources/real-world-evidence/data-analysis-real-world-interrogation-network-darwin-eu> (31 October 2025, date last accessed).
22. OPTIMA. *OPTIMA Consortium*. <https://www.optima-oncology.eu/optima-consortium/> (31 October 2025, date last accessed).
23. UK Parliament. *Constituency data: GPs and GP practices*. <https://commonslibrary.parliament.uk/constituency-data-gps-and-gp-practices/> (31 October 2025, date last accessed).
24. Delmestri A, Prieto-Alhambra D. CPRD GOLD and linked ONS mortality records: Reconciling guidelines. *Int J Med Inform* 2020;**136**:104038. <https://doi.org/10.1016/j.ijmedinf.2019.104038>
25. NHS Digital. *The NHS Research Secure Data Environment Network*. <https://digital.nhs.uk/data-and-information/research-powered-by-data/sde-network> (31 October 2025, date last accessed).
26. European Medicines Agency. *Data Partners onboarded in DARWIN EU (Phase I II III)* [https://www.ema.europa.eu/en/documents/other/darwin-eu-data-partners-onboarded-phases-i-ii-iii\\_en.pdf](https://www.ema.europa.eu/en/documents/other/darwin-eu-data-partners-onboarded-phases-i-ii-iii_en.pdf) (31 October 2025, date last accessed).
27. Observational Health Data Sciences and Informatics (OHDSI). *OMOP Common Data Model (version 5.4) - Observation table*. <https://ohdsi.github.io/CommonDataModel/cdm54.html#observation> (31 October 2025, date last accessed).
28. Union for International Cancer Control (UICC). *TNM Classification of Malignant Tumours*. <https://www.uicc.org/what-we-do/sharing-knowledge/tnm> (31 October 2025, date last accessed).
29. Kahn MG, Callahan TJ, Barnard J *et al*. A harmonized data quality assessment terminology and framework for the secondary use of electronic health record data. *EGEMS (Wash DC)* 2016;**4**:1244. <https://doi.org/10.13063/2327-9214.1244>
30. Strongman H, Williams R, Bhaskaran K. What are the implications of using individual and combined sources of routinely collected data to identify and characterise incident site-specific cancers? a concordance and validation study using linked English electronic health records data. *BMJ Open* 2020;**10**:e037719. <https://doi.org/10.1136/bmjopen-2020-037719>
31. Omar MI, MacLennan S, Ribal MJ *et al*; PIONEER Consortium. Unanswered questions in prostate cancer - findings of an international multi-stakeholder consensus by the PIONEER consortium. *Nat Rev Urol* 2023;**20**:494–501. <https://doi.org/10.1038/s41585-023-00748-9>
32. Tan E, Newby D, Man W *et al*. Comparison of baseline characteristics in patients with prostate cancer who received radical prostatectomy or external beam radiotherapy: target trial versus real-world data population. *Pharmacoepidemiol Drug Saf* 2025;**34**:130. <https://doi.org/10.1002/pds.70186>
33. Whitfield E, White B, Barclay ME *et al*. Differences in recording of cancer diagnosis between datasets in England: a population-based study of linked cancer registration, hospital, and primary care data. *Cancer Epidemiol* 2025;**94**: 102703. <https://doi.org/10.1016/j.canep.2024.102703>
34. Genomics England. *Public OMOP Mappings*. [https://gitlab.com/genomicsengland/genomics\\_england\\_publications/public-omop-mappings](https://gitlab.com/genomicsengland/genomics_england_publications/public-omop-mappings) (31 October 2025, date last accessed).