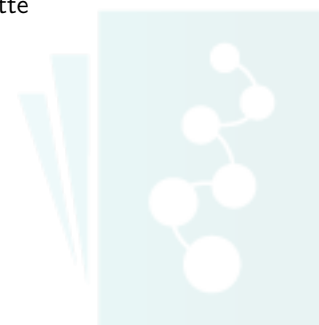


## Supplementary Appendix 1: GUILD report on ASSIGN: The checklist uses data linkage reporting principals from the GUILD Guidance for Information about Linking Data Sets

<b>Data provision</b>		
<b>Concept</b>	<b>Discovery data service (DDS) patient addresses</b>	<b>AddressBase Premium</b>
Population included	<p>Distinct current GP registered patient addresses as at 16<sup>th</sup> November 2020 from 7 CCG GP practices in north east London for persons aged 18 and over.</p> <p>n = 945,196 distinct addresses</p> <p>Reporting on distinct addresses so that the number of patients with the same address does not skew results.</p>	<p>Records for Greater London area plus 8km buffer Epoch 75.</p> <p>n = 10,595,513 (local authority Land and Property Identifier LPI and Royal Mail Delivery Point Address DPA records)</p>
Linkability: how generated	Addresses provided by patients either online or on a paper form when registering with GPs	Master list of addresses sourced from Ordnance Survey, Royal Mail and local authorities
Linkability: how processed	Entered manually by GP practice administrators	Managed and maintained by GeoPlace <sup>2</sup>
Linkability: how quality controlled	Varies by practice: either no quality control, or check against a street list, or Google searches	GeoPlace stringent data quality processes. Run 359 checks on each record before being accepted into the database. BS7666 <sup>3</sup> standard.
Linkability: updates	When informed by patient. Updated addresses are available to Discovery Data Service in real-time	Every 6 weeks
Linkability: cleaning and validation	<p>Address data quality measures calculated. The addresses are reformatted:</p> <ul style="list-style-type: none"> <li>• into eleven standard address object fields: flat, building, number, dependent thoroughfare, street, dependent locality, locality, town, postcode, organisation, vertical</li> <li>• a second version of the eleven standard address object field is created by correcting spelling errors, de-pluralisation, replacing or removing punctuation and lower casing, and removing extraneous words that are unnecessary in the match process, for example, the range of words that are equivalent to the word 'flat' such as 'apartment' or 'maisonette'</li> <li>• positional checking is carried out e.g. the abbreviation 'st' would be mapped to "street" as a spelling correction, but not if it was presented as the first word in a field "St David's" for example would be retained as "St David".</li> </ul> <p>See <a href="https://github.com/endeavourhealth-discovery/uprn-match/tree/master/UPRN/yottadb">https://github.com/endeavourhealth-discovery/uprn-match/tree/master/UPRN/yottadb</a> for address preformatting routines.</p>	<p>The addresses are reformatted:</p> <ul style="list-style-type: none"> <li>• into eleven standard address object fields: flat, building, number, dependent thoroughfare, street, dependent locality, locality, town, postcode, organisation, vertical</li> <li>• the eleven standard address object fields are indexed with single and compound indexes to improve search performance time</li> <li>• the eleven standard address object fields are indexed with performance improving indexes based on semantic equivalence or semantic performance including correcting spelling errors, de-pluralisation, replacing or removing punctuation and lower casing, and removing extraneous words that are unnecessary in the match process, for example, the range of words that are equivalent to the word 'flat' such as 'apartment' or 'maisonette'</li> </ul>
Linkability: replaced with artificial identifiers to reduce disclosure before linkage	N/A	N/A



## Supplementary Appendix 1: Continued

<b>Data linkage</b>		
<b>Concept</b>	<b>DDS patient addresses</b>	<b>AddressBase premium</b>
Process: characteristics used for linkage	Address and postcode	Address and postcode
Process: patterns of missingness	<p>There are 945,196 total <i>distinct</i> addresses of which 804 (0.09%) have a missing or invalid address or postcode<sup>1</sup>.</p> <p><sup>1</sup>An incomplete address &lt;8 characters in length; or contains no alphanumeric characters; or contains the words: unknown, no fixed abode, dummy, nfa, not found, not entitled, overseas, not known, not given, overseas, patient, visitor, unk, address, zz99, @, place of birth, none; or begins with: a special character, london, xx, or x; or does not follow full UK postcode format</p>	N/A
Process: expected range of values after cleaning	N/A	N/A
Process: de-duplication	Duplicate address strings relating to different patient-address pairs removed in previous step. Duplicate addresses that are formatted differently were included because they could not easily be identified as relating to the same address until UPRNs are assigned.	N/A Duplicate versions of UPRN in ABP due to different versions of the same address reflecting aliases and the address life cycle
Process: description of algorithm	<p><b>Reformat</b> Candidate and standard addresses are reformatted as per 'cleaning and validation' section.</p> <p><b>Match</b> Blocking by matching postcode area, potential matching standard addresses are assessed deterministically by applying matching judgement rules in rank order of extent of string manipulation (rank 1 = no manipulation), using a decision tree to determine which string comparison match tests are passed and which fail until all branches are exhausted and the best match is found. These rules mirror human pattern recognition and are coded using e.g. Levenshtein distance<sup>4</sup>, pattern matching (Regex), field swapping and pluralisation. A match is made with one of four overall qualifiers that qualifies the relationship between the candidate address and the matched standard address in relation to approximate geography, or no match is made. The four qualifiers are:</p> <ul style="list-style-type: none"> <li>• Best match: the closest match out of all available</li> <li>• Child: candidate address is a 'child' sub-property of the UPRN it has been matched to</li> <li>• Parent: candidate address is the 'parent' building shell of the UPRN it has been matched to</li> <li>• Sibling: candidate address is a near neighbour of the UPRN it has been matched to</li> </ul> <p><b>Return</b> Where there is a match, the algorithm returns the UPRN, the overall qualifier, the standard address, the match pattern and match rule identifier employed to get that match. The match rule is a label identifying which section of the code made the match, and the match pattern depicts how five address objects were manipulated to achieve the match. These five address objects are merged from the original eleven: flat, building, number, street, postcode. Twelve possible match terms (see Table 1) exist and can be combined in up to 50 different ways on the five address fields. These are restricted to plausible terms, for example, postcodes are never swapped with streets.</p>	

Continued.

## Supplementary Appendix 1: Continued

Data linkage																																																																				
Concept	DDS patient addresses	AddressBase Premium																																																																		
Process: new derived linkage variables	<p>An example of a match pattern is 'Pe,Se,Ne,Bp,Fe'. This means that the postcode, street, number, and flat fields were equivalent matches between the candidate and standard address, and the building field was a partial match between the candidate and standard address. The algorithm is described here: <a href="https://wiki.discoverydataservice.org/index.php?title=UPRN_address_matching_algorithm">https://wiki.discoverydataservice.org/index.php?title=UPRN_address_matching_algorithm</a></p> <p>The algorithm is available for free open-source use here: <a href="https://github.com/endeavourhealth-discovery/ASSIGN">https://github.com/endeavourhealth-discovery/ASSIGN</a></p>																																																																			
Process: blocking methods	By postcode area																																																																			
Record-level indicators of the process	UPRN, qualifier, match rule, match pattern																																																																			
Aggregate linkage results: number of records linked and unlinked	<p>Of 945,196 distinct address strings:</p> <p>924,094 matched (<b>98%</b>)</p> <p>21,102 unmatched (<b>2%</b>)</p>	N/A																																																																		
Aggregate linkage results: comparison of characteristics of linked and unlinked records	<p>Of 924,094 matched, broken down by qualifier:</p> <table border="1" data-bbox="427 931 903 1133"> <thead> <tr> <th>Qualifier</th> <th>Count</th> <th>%</th> </tr> </thead> <tbody> <tr> <td>Best match</td> <td>904,259</td> <td>97.85</td> </tr> <tr> <td>Child</td> <td>9,912</td> <td>1.07</td> </tr> <tr> <td>Parent</td> <td>686</td> <td>0.07</td> </tr> <tr> <td>Sibling</td> <td>9,237</td> <td>1.00</td> </tr> <tr> <td>Total matched</td> <td>924,094</td> <td></td> </tr> </tbody> </table> <p>Address characteristics:</p> <table border="1" data-bbox="427 1171 1477 1570"> <thead> <tr> <th>Characteristic</th> <th>Total</th> <th>Linked</th> <th>Unlinked</th> </tr> </thead> <tbody> <tr> <td>Total</td> <td>1,549,669</td> <td>1,425,497</td> <td>124,172</td> </tr> <tr> <td>Of which:</td> <td></td> <td></td> <td></td> </tr> <tr> <td>E postcode %</td> <td>61.2</td> <td>61.2</td> <td>62.0</td> </tr> <tr> <td>N postcode %</td> <td>7.3</td> <td>7.3</td> <td>9.0</td> </tr> <tr> <td>R postcode %</td> <td>18.7</td> <td>19.0</td> <td>6.0</td> </tr> <tr> <td>I postcode %</td> <td>12.3</td> <td>12.3</td> <td>11.8</td> </tr> <tr> <td>Other postcode %</td> <td>0.5</td> <td>0.3</td> <td>8.6</td> </tr> <tr> <td>Address begins with numeric character %</td> <td>75.9</td> <td>76.5</td> <td>52.7</td> </tr> <tr> <td>Address begins with alphabetic character %</td> <td>24.0</td> <td>23.5</td> <td>46.6</td> </tr> <tr> <td>Address begins with special character %</td> <td>0.0</td> <td>0.0</td> <td>0.7</td> </tr> <tr> <td>Invalid address or postcode %</td> <td>0.1</td> <td>0.0</td> <td>3.5</td> </tr> </tbody> </table> <p>There are higher proportions of 'Other' postcodes, addresses beginning with an alphabetic character (i.e. a flat rather than a house) or a special character, and invalid addresses or postcodes in unmatched compared to matched. Differences between matched and unmatched addresses across all characteristics were found to be significant using chi square tests, but this could be attributable to the large sample size. Patient and registration characteristics are compared in section 'Population characteristics' of the paper.</p>		Qualifier	Count	%	Best match	904,259	97.85	Child	9,912	1.07	Parent	686	0.07	Sibling	9,237	1.00	Total matched	924,094		Characteristic	Total	Linked	Unlinked	Total	1,549,669	1,425,497	124,172	Of which:				E postcode %	61.2	61.2	62.0	N postcode %	7.3	7.3	9.0	R postcode %	18.7	19.0	6.0	I postcode %	12.3	12.3	11.8	Other postcode %	0.5	0.3	8.6	Address begins with numeric character %	75.9	76.5	52.7	Address begins with alphabetic character %	24.0	23.5	46.6	Address begins with special character %	0.0	0.0	0.7	Invalid address or postcode %	0.1	0.0	3.5
Qualifier	Count	%																																																																		
Best match	904,259	97.85																																																																		
Child	9,912	1.07																																																																		
Parent	686	0.07																																																																		
Sibling	9,237	1.00																																																																		
Total matched	924,094																																																																			
Characteristic	Total	Linked	Unlinked																																																																	
Total	1,549,669	1,425,497	124,172																																																																	
Of which:																																																																				
E postcode %	61.2	61.2	62.0																																																																	
N postcode %	7.3	7.3	9.0																																																																	
R postcode %	18.7	19.0	6.0																																																																	
I postcode %	12.3	12.3	11.8																																																																	
Other postcode %	0.5	0.3	8.6																																																																	
Address begins with numeric character %	75.9	76.5	52.7																																																																	
Address begins with alphabetic character %	24.0	23.5	46.6																																																																	
Address begins with special character %	0.0	0.0	0.7																																																																	
Invalid address or postcode %	0.1	0.0	3.5																																																																	
Aggregate linkage results: representativeness of the linked data set	See paper section 'Bias in UPRN match success'																																																																			
Aggregate linkage results: flow diagram of linkage steps	N/A – the linkage steps pathway is different for different addresses depending on the content and required manipulation of the address string																																																																			

Continued.

Supplementary Appendix 1: Continued

**Data linkage**

Concept	DDS patient addresses	AddressBase premium												
Linkage accuracy: how error rates were estimated	Algorithm applied to two ‘gold-standard’ external reference data sets. 1) 9,177 Welsh local authority addresses. 2) 9,475 Tower Hamlets local authority addresses  True false positive matches, false matches, missed matches, and true negative matches are quantified to calculate: <ul style="list-style-type: none"> <li>• Positive Predictive Value (PPV) or Precision - the proportion of record pairs classified by the algorithm as links that are true matches</li> <li>• Sensitivity or Recall– the proportion of true matches that are correctly classified as links.</li> <li>• The F-measure – The harmonic mean between positive predictive value and sensitivity. Often used to compare the overall efficiency of a method</li> </ul>													
Linkage accuracy: estimates of error rates	<table border="1"> <thead> <tr> <th>Measure</th> <th>DDS address linkage results on Welsh gold-standard addresses</th> <th>DDS address linkage results on Tower Hamlets gold-standard addresses</th> </tr> </thead> <tbody> <tr> <td>Sensitivity</td> <td>0.999</td> <td>0.999</td> </tr> <tr> <td>PPV</td> <td>0.996</td> <td>0.998</td> </tr> <tr> <td>F-measure</td> <td>0.997</td> <td>0.998</td> </tr> </tbody> </table>		Measure	DDS address linkage results on Welsh gold-standard addresses	DDS address linkage results on Tower Hamlets gold-standard addresses	Sensitivity	0.999	0.999	PPV	0.996	0.998	F-measure	0.997	0.998
Measure	DDS address linkage results on Welsh gold-standard addresses	DDS address linkage results on Tower Hamlets gold-standard addresses												
Sensitivity	0.999	0.999												
PPV	0.996	0.998												
F-measure	0.997	0.998												
Disclosure controls	Addresses and UPRNs remain in the identifiable zone of Discovery Data Service only. UPRNs are pseudonymised into Residential Anonymous Linking Fields for third party use													

<sup>1</sup>Gilbert, R., Lafferty, R., Hagger-Johnson, G., Harron, K., Zhang, L.C., Smith, P., Dibben, C. and Goldstein, H., 2017. GUILD: GUIDance for Information about Linking Data sets. *Journal of Public Health*, 2017 Mar 28:1–8.

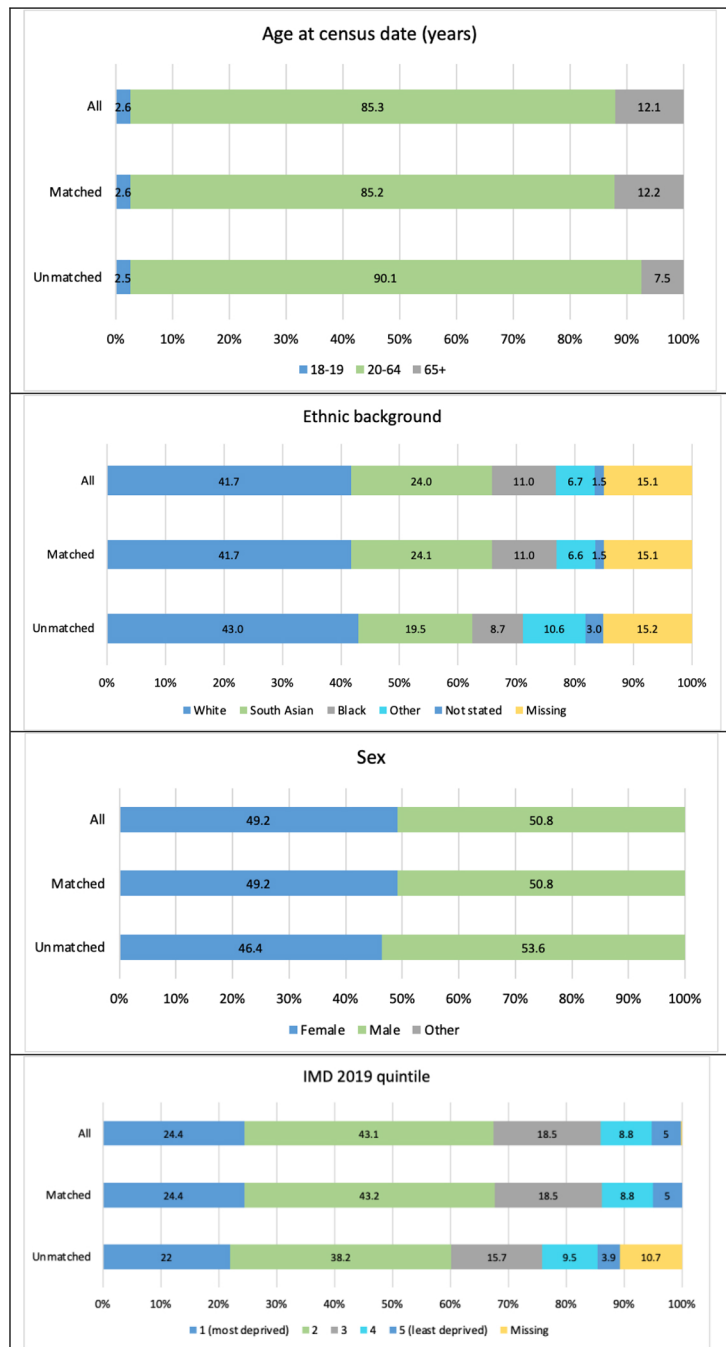
<sup>2</sup>[www.geoplace.co.uk](http://www.geoplace.co.uk)

<sup>3</sup><https://www.aligned-assets.co.uk/british-standard-bs7666/>

<sup>4</sup>[https://en.wikipedia.org/wiki/Levenshtein\\_distance](https://en.wikipedia.org/wiki/Levenshtein_distance)



Supplementary Appendix 2: Summary characteristics of the study population according to whether patient address was matched or not matched to a UPRN by the ASSIGN algorithm



UPRN = Unique Property Reference Number.  
 IMD = Index of Multiple Deprivation.



Supplementary Appendix 3: UPRN match rates and absolute differences in proportion matched with respect to reference category for all explanatory variables N = 1,757,018

	Number <i>n</i>	Address-matched to UPRN (%)	Absolute difference relative to reference group (%)
<b>Age at census date 16/11/2020 (years)</b>			
Missing	8,116	99.62	0.06
<b>&gt;1</b>	<b>50,740</b>	<b>99.56</b>	<b>Ref</b>
1–14	133,371	99.33	–0.22
15–29	570,251	98.06	–1.49
30–64	929,452	98.71	–0.85
65–84	59,973	98.77	–0.85
85 and over	<b>5,115</b>	<b>96.72</b>	<b>–2.84</b>
<b>Ethnic background</b>			
Missing	265,524	98.56	–0.08
<b>British</b>	<b>382,170</b>	<b>98.64</b>	<b>Ref</b>
African	100,743	98.68	0.03
Any other Asian background	61,521	98.38	–0.27
Any other Black background	44,131	99.01	0.37
Any other White background	337,905	98.4	–0.24
Any other ethnic group	52,823	98.42	–0.22
Any other mixed background	15,018	97.88	–0.77
Bangladeshi	145,920	99.28	0.64
Caribbean	48,203	99.16	0.51
Chinese	<b>21,961</b>	<b>95.51</b>	<b>–3.14</b>
Indian	121,134	98.51	–0.13
Irish	13,113	98.41	–0.24
Not stated	<b>26,196</b>	<b>97.09</b>	<b>–1.56</b>
Pakistani	93,538	98.9	0.25
White and Asian	4,947	98.08	–0.56
White and Black African	9,971	97.9	–0.74
White and Black Caribbean	12,200	98.21	–0.43
<b>Sex</b>			
<b>Female</b>	<b>864,337</b>	<b>98.65</b>	<b>Ref</b>
Male	892,638	98.49	–0.16
Other	43	95.35	–3.3
<b>IMD 2019 quintile</b>			
Missing	<b>3,502</b>	<b>23.5</b>	<b>–75.21</b>
<b>1 (most deprived)</b>	<b>428,373</b>	<b>98.71</b>	<b>Ref</b>
2	757,212	98.74	0.02
3	325,075	98.79	0.08
4	154,523	98.45	–0.26
5 (least deprived)	88,333	98.88	0.17
<b>GP registration duration (quartiles)</b>			
Missing	<b>8,116</b>	<b>99.58</b>	<b>1.94</b>
<b>1 (shortest)</b>	<b>437,228</b>	<b>97.64</b>	<b>Ref</b>
2	437,422	98.36	0.72
3	<b>437,603</b>	<b>98.92</b>	<b>1.28</b>
4 (longest)	<b>436,649</b>	<b>99.36</b>	<b>1.72</b>

Continued.

## Supplementary Appendix 3: Continued

	Number <i>n</i>	Address-matched to UPRN (%)	Absolute difference relative to reference group (%)
<b>Number of GP registrations in preceding 12 months</b>			
1	<b>1,595,729</b>	<b>98.58</b>	<b>Ref</b>
2	144,755	98.61	0.03
3 or more	16,534	97.67	-0.91
<b>Number of address changes in preceding 12 months</b>			
1	<b>1,316,956</b>	<b>98.98</b>	<b>Ref</b>
2	<b>343,808</b>	<b>97.89</b>	<b>-1.09</b>
3 or more	<b>96,254</b>	<b>95.41</b>	<b>-3.57</b>
<b>GP system</b>			
Missing	4,960	99.62	0.83
<b>EMIS</b>	<b>1,629,199</b>	<b>98.79</b>	<b>Ref</b>
SystemOne	<b>87,783</b>	<b>94.39</b>	<b>-4.4</b>
Vision	35,076	98.86	0.08
<b>Clinical Commissioning Group</b>			
<b>Newham</b>	<b>326,386</b>	<b>99.16</b>	<b>Ref</b>
Barking & Dagenham	168,008	98.59	-0.57
City & Hackney	259,973	98.25	-0.91
Havering	221,328	99.38	0.22
Redbridge	251,128	98.61	-0.55
Tower Hamlets	<b>278,520</b>	<b>97.7</b>	<b>-1.46</b>
Waltham Forest	251,675	98.35	-0.81

Quartile definitions for GP registration duration: Quartile 1 (shortest): 0–32 months; Quartile 2: 33–77 months; Quartile 3: 78–183 months; Quartile 4 (longest) > 184 months.

EMIS: Egton Medical Information Systems.

Reference groups and values with an absolute match rate difference to the reference group of >1% are in bold.

