

# Technical Guide: Estimates of children with a parent in prison

## 1. Introduction

These estimates of children under 18 with a parent in prison are designated as [Official Statistics in Development](#) and were produced in accordance with the [Code of Practice for Official Statistics](#).

This document provides a comprehensive guide to the statistics. It covers:

- An explanation of the data sources and quality used to produce the estimates;
- The methodology adopted to compile the estimates (including data linking and Natural Language Processing);
- Assumptions and limitations of the data and analysis; and
- Users of the statistics.

## 2. Data sources and quality

### 2.1 Summary of data sources

These statistics draw from data from administrative databases across His Majesty's Prisons and Probation service (HMPPS). The five specific data sources used were:

#### Structured data sources

1. The Prison National Offender Management Information System (p-NOMIS) **contact lists**
2. The Offender Assessment System (OASys) **Basic Custody Screening Tool (BCST), Part 1**
3. The probation case management system, national Delius (nDelius) **personal circumstances flags**

#### Unstructured data sources

4. The Prison National Offender Management Information System (p-NOMIS) **prison case notes**
5. The probation case management system, national Delius (nDelius) **probation case notes**

In addition, we carried out a data-matching pilot with His Majesty's Revenue and Customs (HMRC) using the following data:

6. HMRC Child Benefit records
7. HMRC Working Tax Credit records – explored but not used to create final estimates

p-Nomis holds information on prisoners, their movements, and their activity while in public prisons. This database was used for identification of the prisoner cohort and prisoner contact lists.

OASys contains assessments for offenders at any point during their journey through the Criminal Justice System, including while in custody or on probation. This includes the Basic Custody Screening Tool (BCST).

nDelius holds personal circumstances flags for a range of topics. For this work, the relevant flag category is “Has Dependents”.

The quality of the information held in these operational databases, while generally high, does differ by data field, being dependent on its frequency of use and importance in the day-to-day running of the prisons and probations systems. Specific detail on the quality of each data field that is used has been included in the following section.

HMPPS data has been linked to HMRC Child Benefit data for analysis. The quality of data drawn from all the above sources was used to produce aggregated findings only. Analysts have not directly explored the records or circumstances of individual prisoners.

## **2.2 Use of administrative data**

The use of administrative operational data allows the possibility of information to be included on all prisoners, rather than a sample. In addition, the data has already been collected for operational purposes and so does not require additional resource. Details of all administrative data sources used in the production of this release can be found in the MoJ Statement of Administrative Sources.<sup>1</sup>

However, this data has the same limitations as any other data that comes from large-scale administrative record systems, meaning that there may be mistakes with entering and processing the data. Also, in terms of its direct relevance to estimating how many prisoners have children, it is important to note that the way dependants are defined may vary across the administrative data sources used, as well as how operational staff who collect the information interpret these definitions.

## **2.3 Specific datasets used**

### **2.3.1 p-NOMIS – contact lists**

p-NOMIS holds details of prisoners’ contacts who come to visit them in prison. When this is recorded on the system, the name, date of birth and relationship type can be included on the record. We have filtered the contact lists to only include sons, daughters, stepsons and stepdaughters. As we have defined a child as anyone aged under 18 on the last day of the cohort period (i.e. as of 1 October 2022), we have excluded children aged 18 and over. The quality of the prisoner contacts information is generally assumed to be good. However, if a prisoner is not visited by a family member, this information will not be recorded. Data drawn from the contact lists has been used to produce aggregate findings only.

---

<sup>1</sup> Ministry of Justice (2016), [Statement of Administrative Sources](#) (PDF, 624KB, 58 pages)

### **2.3.2 p-NOMIS – case notes**

Case notes are a section of p-NOMIS where staff can record information about an offender's behaviour, progress, and other relevant details.<sup>2</sup> Case notes provide a comprehensive and up-to-date overview of an offender's situation.

All staff who have contact with an offender and who have access to p-NOMIS must update case notes on a regular basis. Management checks are in place to ensure frequency and quality of entries in case notes.

### **2.3.3 OASys – Basic Custody Screening Tool (BCST) Part 1**

Part 1 of the BCST is conducted within 72 hours of arrival into custody, whether the prisoner is remanded into custody or has received a custodial sentence. The BCST is designed to promptly identify key issues and needs. The relationships section of the BCST captures information on an individual's parental/caregiving status for those aged under 18 years. There are known limitations of the BCST coverage including a prisoners' (un)willingness to disclose information and (lack of) time to complete the BCST by assessors. Questions and assessment processes conducted with offenders are under continuous review by the MoJ to reflect learning from operational colleagues and people with lived experience of prison.

Many offenders do disclose whether they have children at the point when the BCST is conducted and so, despite some data quality issues, the BCST is often cited as the main source of understanding which prisoners have children. Including this data allows us to identify additional offenders who are parents. Examining the overlap with other data sources helped us understand the completeness of information about dependants across different records.

### **2.3.4 nDelius – personal circumstance flags**

Flags can be added to an offender's probation record on nDelius. They are to store additional relevant information about an offender. The primary purpose of these flags is for officers to be aware of the circumstances of the offender and how this may affect the offender's ability to comply with their probation requirements.

The category used in this work was the "Has dependents" flag. When the flag is added, information on the names, ages and circumstances of the dependants can be added in the comments field. However, the guidance for officers entering data does not specify that this information should only include dependants under 18, and so the flag could relate to caring for a dependent adult rather than a child. From initial exploration of the free-text fields and engagement with stakeholders, we expect that a very small number of these records relate to adults rather than children.

### **2.3.5 nDelius – case notes**

Case notes should be a record of contact between the probation practitioner and the person on probation, as well as other contacts the practitioner may have with other

---

<sup>2</sup> NOMS Agency Board (2014), [p-NOMIS instructions](#)

parties.<sup>3</sup> Contact should be recorded within one working day. Case notes are expected to be an explicit record of the nature, location and time of any contact. Records should distinguish between fact and opinion and contain sufficient information to support probation practitioner tasks.

### **2.3.6 HMRC Child Benefit data**

Child Benefit is provided to families responsible for bringing up a child under 16, or under 20 if they are in approved education or training.<sup>4</sup> Only one person can claim Child Benefit for a child. There is no limit to how many children you can claim for. Because age is part of the qualification criteria, benefit records also contain the numbers and ages of children as well as the address.

As of August 2023, 95.6% of all child benefit records related to children aged under 18. As the data therefore relates mostly to children, it should be a reliable indicator of whether an individual has dependants under 18.<sup>5</sup>

We note that Child Benefit is claimed by the legal guardian so in some cases the claimant may not be the parent.<sup>6</sup> Although most claimants are parents, they could also be grandparents, guardians and others if the child has been removed from parental care. Child Benefit take-up rates are now at 90% (as of May 2023).<sup>7</sup>

A data-linking exercise was previously carried out matching Department for Work and Pensions (DWP) Child Benefit data to female offenders in Police National Computer (PNC) data to get an estimate of those with child dependants in 2012.<sup>8</sup> Between 24% and 31% of all female offenders were estimated to have one or more child dependants. The number of children associated with each claimant was identified in the benefit records with a mean value of 1.9. This exercise therefore suggested that a similar matching exercise linking information on offenders held in HMPPS data to HMRC Child Benefit data would be beneficial.

### **2.3.7 HMRC Working Tax Credit data**

Working Tax Credit generally requires all members of a household to provide information in order to assess eligibility.<sup>9</sup> It was therefore identified as a good option to match both to the (typically) male offender and the (typically) female partner. It also has a Child Tax Credit component so contains details of children in the

---

<sup>3</sup> HM Prison & Probation Service, [Probation Service – National Standards 2021 – guidance update June 2022](#) (MS Word document, 4.36MB)

<sup>4</sup> See guidance provided on GOV.UK: <https://www.gov.uk/child-benefit>.

<sup>5</sup> HM Revenue and Customs, released April 2024, GOV.UK, [Child Benefit Statistics: annual release, data at August 2023](#)

<sup>6</sup> You can claim Child Benefit if you're responsible for a child under 16 (or under 20 if they're in approved education or training)" as referred to on this factsheet: [https://assets.publishing.service.gov.uk/media/5cbf1b7d40f0b63cacd6dcaf/Child\\_Benefit\\_factsheet.pdf](https://assets.publishing.service.gov.uk/media/5cbf1b7d40f0b63cacd6dcaf/Child_Benefit_factsheet.pdf)

<sup>7</sup> HM Revenue and Customs, released April 2024, GOV.UK, [Child Benefit Statistics: annual release, data at August 2023](#)

<sup>8</sup> Ministry of Justice, released October 2015, GOV.UK, [Female offenders and child dependents](#)

<sup>9</sup> See guidance provided on GOV.UK: <https://www.gov.uk/working-tax-credit>.

household. Linking to Working Tax Credit (WTC) was also explored as an option for linking to HMPPS data. However, initial results were poor given that many of families during the study time period were transitioning over to Universal Credit and hence were missing from the WTC data.<sup>10</sup> This did not affect Child Benefit which is a stand-alone benefit.

## 2.4 Removed records

Records with conflicting information in characteristics for the same individual due to data quality issues have been omitted from the counts in the tables in both the main report and Technical Guide. Individuals may have conflicting records listed for multiple reasons including clerical error, errors when matching individuals across data sources, or changing circumstances (i.e. sentence length may change after an appeal). This affects 6% of the prisoner cohort.

## 3. Data governance

The BOLD programme has established procedures for the effective governance of data it uses across all pilot projects. You can find [BOLD's Privacy Notice on GOV.UK](#).

### 3.1 Governance

Analysis and research using data collected in operational systems across HMPPS is covered by pre-existing Data Protection Impact Assessments (DPIA). This work is covered by these under the research and analysis purpose. Statistics are in aggregate form only for the purposes of understanding offenders; information about any specific individual is not of interest.

### 3.2 Confidentiality

This statement sets out the arrangements in place for protecting persons' confidential data when statistics are published or otherwise released into the public domain. The [Code of Practice for Statistics](#) states that:

*“Organisations should look after people’s information securely and manage data in ways that are consistent with relevant legislation and serve the public good.”<sup>11</sup>*

To comply with this and with the Data Protection Act of 2018 and to maintain the trust and co-operation of those who use these Official Statistics in Development, the following provisions have been put in place:

- Private information collected by MoJ is stored in line with our data security policies.
- Electronic data is held on password-protected networks.

---

<sup>10</sup> See the National Audit Office's (NAO's) February 2024 report [Progress in implementing Universal Credit](#)

<sup>11</sup> Statistics Authority (2018, updated May 2022), [Code of Practice for Statistics edition 2.1 – T6: Data governance](#)

- All new staff undergo security vetting before receiving access to data systems and all staff undertake mandatory training on information responsibility annually.

Some counts may have been removed for Statistical Disclosure Control purposes. In line with MoJ and GSS guidance, assessment of the risk of disclosure considers the following:

- Level of aggregation (including geographic level) of the data;
- Size of the population;
- Likelihood of an attempt to identify; and
- Consequences of disclosure.

### **3.3 Engaging the public**

Public trust around how data is shared is critical for BOLD, and we partnered with the Centre for Data Ethics & Innovation (CDEI), and the research company Britain Thinks, to undertake extensive engagement with affected groups, trusted intermediaries, and the general public. The results of this exercise, and what we have learnt from listening to the public, have tangibly informed the design of the BOLD programme and has been [published by the CDEI](#).

## **4. Methodology**

The methods used to develop the statistics in this report are complex and involve probabilistic linkage (section 4.1), the process by which personal records from one data set are attached to personal records from another. Extraction of information from case notes involved Natural Language Processing and inference (section 4.2). The HMRC pilot study is explained in more detail in section 4.3.

### **4.1 Data-linking**

The three main HMPPS databases (p-NOMIS, nDelius, and OASys) set out in section 2 do not have a common personal identifier to enable the same individual to be identified across the systems. This means the data cannot be linked in a straightforward way. In addition, multiple records may be associated with the same individual in one database. There is no unique identifier available to reliably link records for the same person from within and between the databases.

Splink has therefore been used to identify unique, deduplicated offender records and link the data across datasets. Splink is based on probabilistic linkage and was developed internally by MoJ's data linking team. Further details about Splink are available in [Data First: An Introductory User Guide](#) (PDF, 951KB, 36 pages) and [Data First: Criminal Courts Linked Data](#).

### **4.2 Free-text analysis**

As this work contains some novel techniques, we have described the approach used in full. The methods used to extract information from the unstructured (free-text) prison and probation officer case notes involved Natural Language Processing (NLP)

techniques and use of Large Language Models (LLMs). A useful guide on LLMs and their relevance for statistics has been provided by Eurostat.<sup>12</sup>

In summary, given the case notes of an individual, we used an NLP tool to select sentences which contain mentions to children using a list of keywords (see section 4.2.1). The sentences containing child mentions were processed using a Natural Language Inference (NLI) model (section 4.2.2) which enables us to determine an indication of parental status.

To arrive at a binary prediction on whether a given sentence suggests the prisoner has children, the output from the model (a likelihood the prisoner has a child, given the sentence) is compared to a **threshold**. A probability greater than the threshold from any sentence results in the offender being marked positive as having children.

The NLI model has been fine-tuned (section 4.2.5) to improve its effectiveness based on synthetic text data generated by an LLM (section 4.2.4), which was in turn prompted using indicative case note data (section 4.2.3). The indicative data was produced by HMPPS operational staff.

The final free-text model (as referred to in the main report), was validated against labelled data which involved manual classification of sentences to state whether the individual referred to is a parent or not (section 4.2.6).

#### 4.2.1 Child mention extraction

We used the Natural Language Processing information extraction tool [spaCy matcher](#) to identify sentences with mentions to children from case notes based around pre-defined keywords: son, daughter, child, kid, stepson, stepdaughter, stepchild. However, not all sentences containing these keywords mean that the offender has children. For example, a sentence may refer to the offender's own childhood or their childish behaviour. We used a Natural Language Inference model in an attempt to verify the nature of the keyword mention. System-generated notes were excluded from this analysis in order to reduce the amount of data that needed to be processed.

#### 4.2.2 Natural Language Inference

Natural Language Inference (NLI) is the task of determining whether a hypothesis is true, false, or undetermined. For example, given the hypothesis "this person has children", then the sentence "A stated he has two children" should be true, "B has a partner but no children" false, and "C behaved like a child" undetermined.

The specific model used was the [DeBERTa v3 small model](#), however this was unable to correctly pick up on sentences which contained language specific to the prison and probation systems and the domain of family and children. In order to improve NLI performance on our data, we further fine-tuned our model on a

---

<sup>12</sup> See Eurostat, Buono D, Felecan M and Tessitore C (2024), [An introduction to Large Language Models and their relevance for statistical offices – 2024 edition](#)

synthetically generated dataset using indicative case notes. using indicative case notes.

#### 4.2.3 Indicative case notes

This project was completed in collaboration with prison and probation staff who produced 20 examples of case notes each around a page long, that were entirely fictional, bearing no relation to real individuals. We refer to these as *indicative* case notes. These indicative case notes were written in the style and structure typically used and include probation and/or prison officer terminology, abbreviations, different formats (paragraphs, spaces, headings, indexes, labels), symbols, slang and spelling mistakes.

The officers were asked to make reference to the offender having children in some of their notes and not in others. Sentences from these indicative notes were used to design prompts for synthetic data generation.

#### 4.2.4 Synthetic data generation

We used the LLM [Mistral 7B Instruct](#) to generate synthetic sentences for fine-tuning our model in order to make it more accurate at identifying parents in the free text. Portions of the indicative case notes were used to prompt the model to generate sentences in the same format in the following categories:

- Implies the offender has children.
- Implies the offender does not have children.
- Contains children keywords but no information on whether they have children.

An example for the category **implying the offender has children** is below:

**Prompt -**

*“Below are some examples of indicative case notes written by prison officers:*

*Emma and Max Cranmer are excited / anxious but it is positive that they are investing in themselves and exploring childcare options.*

*Kai Harris would like to feel listened to and included in decisions about his son, Chris and unborn baby with Rose.*

*I reassured her that she had done the right thing by making sure that her kids had somewhere safe to stay and that this is a really positive step.*

*Generate an example sentence from a case note which implies the prisoner has a child or stepchild under 18. The prisoner's name is David STOKES”*

**Model output -**

*“In my conversation with David, he expressed concern about arrangements for his daughter's visits during her school holidays.”*

Relevant sentences were randomly sampled from the indicative case notes (and some additional sentences written by MoJ data scientists) to increase the variation in



generated notes. Tables of popular names were gathered from the Office for National Statistics and three were injected into the prompt on each run.<sup>13</sup> In some cases the prisoner's name was replaced with initials to reflect the style seen used by probation officers.

We generated approximately 10,000 sentences overall, and the resulting synthetic data was used to fine tune the model. The approach provides a much more cost-effective way to curate training data than traditional clerical labelling. This approach produced 10,000 labelled sentences, and in turn an optimised fine-tuned model, in only a few days. It is estimated that to produce 10,000 labelled sentences split amongst a pool of 10 annotators would likely take several months to achieve ordinarily.

#### **4.2.5 Model fine-tuning**

The free-text model was trained with a synthetic dataset of 10,000 sentences generated by an LLM. As such, there may be unaccounted differences in distribution between the real and synthetic datasets, potentially resulting in poor model fit (overfitting and generalisation errors); however, we have evaluated and deemed the model's performance to be good (refer to section 4.2.6). Further development of this work can gain a better understanding of the discrepancy by explicitly covering more extreme/edge cases in the evaluation.

The choice of the threshold can result in large variations in the aggregated count of prisoners with children. The predicted number of prisoners aged 35 and under with children with a high threshold of 0.99 suggested 60% of prisoners with children, while a lower threshold of 0.95 suggested 76% of prisoners with children. Further work would be required to understand the model sensitivity and to fully assess and set the optimal threshold. This would build on existing quality tests on the current model which underwent five states of fine-tuning before satisfactory results were achieved.

#### **4.2.6 Model evaluation**

The final free-text model referred to in the main report (sections 3 and 4) was tested on a different set of real data that was labelled by human experts. This was created by selecting a sample of offenders and displaying text chunks (one or two sentences) to annotators who marked the chunks as having children or not or not containing any mention of children. The sentences were sampled based on the age and sex of the prisoner to account for differences in language used and circumstances (if any). Overall, this resulted in ~1,100 labelled sentences which were used to measure the NLI performance.

The model correctly identified parents 95% of the time (precision) and covered 40% of sentences which imply a child (recall). We have opted for prioritising precision

---

<sup>13</sup> Names were gathered from the Office for National Statistics' list of popular baby names in England and Wales for 2021: Official for National Statistics (ONS), released March 2023, ONS website, [.xlsx, Baby names in England and Wales: 2021](#)

over recall as there are many case notes for each prisoner and if any one of them implies a child, we categorise the prisoner as having a child. We expect that this process improves the recall significantly, as there will be sentences where the implication is more obvious.

### **4.3 Pilot linkage HMPPS data to HMRC benefits data**

#### **4.3.1 HMRC Child Benefit data**

The child benefit data was taken from November of each year from 2017 to 2022. There were 7.8 million unique child benefit claimant records matched against offender data and offenders' partners' data. Each child benefit claimant dataset contains the following columns: surname, forename, date of birth, postcode and National Insurance number.

This sample from the prisoner cohort was randomly selected. However, given the majority of the prison population is male (96%),<sup>14</sup> this sample was weighted to include more females than in the standard prison population so as to obtain sufficient information about how well female offenders, as well as male offenders, match to Child Benefit (150 female prisoners:850 male prisoners). Additionally, the sample included the prisoner's partner details when listed in order to include more females for matching. Around a third of prisoners in the sample had partner details listed. Prisoner ages were broadly representative of the prison population.

Probabilistic matching was undertaken in R using the [fastLink function](#) based on two criteria:

- Criteria A: offender matches based on forename, surname, date of birth and postcode.
- Criteria B: offender partner matches using partner forename, partner surname and partner date of birth.

The results of the HMRC data matching pilot can be found in the main report. We provide further information here on summary of linkage by sex (table A1). Where offender parents have successfully been linked to a Child Benefit record, it is typically the female that matches. The majority of matches for male prisoners were based on their female partners while the majority of matches for female prisoners were based on them rather than their male partners. Overall, 94% of all matches to Child Benefit were based on either matching the female offender or female partner (noting however that females were oversampled as described above). Given that the success of the matching is based on the female parent and demographics of the prison population are majority male, successfully matching to the entire prisoner cohort would therefore likely depend on the presence of female partner information within prisoner records.

---

<sup>14</sup> See MoJ, released November 2022, GOV.UK, [Women and the Criminal Justice System 2021](#)

**Table A1: Results of a matched sample of (unique) HMPPS offenders and partners to HMRC Child Benefit records between 2017 and 2022, broken down by sex**

<b>Number of matches</b>	<b>Count</b>	<b>Percent</b>
Number of matches based on female offender	25	14%
Number of matches based on male offender	8	4%
Number of matches based on female partner	142	80%
Number of matches based on male partner	4	2%
Number of sample records matched to Child Benefit <sup>1</sup>	<b>177</b>	<b>100%</b>

<sup>1</sup>Total matches do not sum to 177 as some prisoners matched on both offender and partner so are included twice.

### **4.3.2 Adjustment for undercount**

The first stage of this work involved a direct count of prisoners with children in the structured data fields as well as from free-text contact notes. Given the differences in prisoners counted between the different structured data fields, we recognised there was a possibility of undercount. Moreover, the free-text model applied to adults aged 35 and under so there was known undercount for older adults who will also have children aged under 18.

In order to adjust for undercount, an estimate of prisoners with children was produced through a series of steps which accounted for information from the HMRC data-matching exercise and the extrapolation of the free-text model results to older adults.

We acknowledge there are limitations with this approach. In particular, we assumed the basis for the ratio used to inflate the structured data count was independent of age, i.e. that specific age groups are no more likely to disclose children in case notes than they are in the structured data. Counts for female offenders are small so we have also had to assume the ratio is independent of sex. We note that a direct count of all prisoners with children is preferred and any adjustment for undercount should be minimal.

The series of steps were as follows:

Step 1. Calculated the ratio in additional parents identified in the free-text structured data which is an additional 124% of the parents identified in the structured data. See Table A2.

Step 2. Applied the resulting ratio to inflate the parents identified in the structured data fields for age groups aged 36 and over.

Step 3. An additional 27,967 parents were calculated based on the above.

Step 4. For every 1,000 prisoners we estimate that 41 parents are missed based on the HMRC pilot. This is an additional 5,723 parents to the 139,592 prisoner cohort.

Step 5. This provides an overall adjusted estimate of prisoner parents as 108,990.

**Table A2: Estimated number of prisoners aged 35 and under with dependants by data source type**

Source	Count of prisoners with dependants
Structured data fields	22,525
Free-text case notes	48,902
Increase from free text	27,967
<i>Ratio (for every 1 prisoner with children under 18 in the structured data, there are 2.24 in the free text)</i>	1: 2.24

Table A3 indicates that for females, just under two thirds (63%) could be identified from the structured data sources, and just under two thirds (63%) of female parents were identified in the text data sources (noting that text data sources only apply to ages 35 and under).

**Table A3: Count of parents identified in each data source broken down by sex; percentages are calculated as a proportion of the total number of parents identified for that sex**

Data Source	Female	Female %	Male	Male %
Contact lists	808	19%	23,821	34%
BCST	1,900	45%	23,523	34%
nDelius flag	680	16%	10,959	16%
<b>Structured total</b>	<b>2,669</b>	<b>63%</b>	<b>43,639</b>	<b>62%</b>
Nomis text	1,946	46%	28,648	41%
nDelius text	2,482	59%	44,419	63%
<b>Text total</b>	<b>2,668</b>	<b>63%</b>	<b>46,234</b>	<b>66%</b>
<b>Total<sup>2</sup></b>	<b>4,208</b>	<b>100%</b>	<b>70,067</b>	<b>100%</b>

<sup>2</sup>Note: totals may not add up because individuals may be identified in multiple sources

For males, the pattern was very similar: up to two thirds (62%) of males could be identified as parents from the structured data sources, and two thirds were identified in the text data sources (66%).

However, female parents were more likely than male parents (45% versus 34%) to be picked up via the BCST. Male parents were more likely than female parents (34% vs 19%) to be identified in the contact lists.

## 5. Future development of these statistics

As this is an Official Statistics in Development publication, these statistics are in their testing phase. We are therefore consulting on the methodology used for future iterations of these statistics. As a general overarching goal, the main improvements to the statistics should come from a better linked dataset that identifies the link between parents and children. However, in the absence of this dataset, further work on these statistics could seek to minimise sources of undercount and maximise

effectiveness of the counting methodology via improvements to the free-text model. Development could take the following routes and we welcome feedback on any of these:

<b>Approach</b>	<b>Advantages</b>	<b>Limitations</b>
1. Development of the free-text model to ascertain age of dependants.	Will provide a more accurate direct count of prisoners with children as the free-text model can be applied to ALL adults not just those aged 35 and under.	Age of child may not always be captured in the free text and/or structured data fields.
2. Run the free-text model for all ages and apply statistical weighting by age group according to expected proportions of adults with children under 18.	May be more reliable as would remove our current assumption around ratio remaining constant across age groups.	Requires linking to Other Government Data (OGD) data in order to validate estimates and to extract the data on children.
3. Provide a range and run free-text model with both (i) age 35 and under restriction and (ii) all offender ages to provide lower and upper bounds with estimate of prisoners with children sitting somewhere in between.	Will better quantify levels of uncertainty around the estimate.	Requires linking to OGD data in order to extract the data on children.
4. Further development of the HMRC matching exercise (i) analysing the match-rate based on structured data <u>and</u> free-text (ii) matching to the full prisoner cohort.	Will provide a direct count of children with a parent in prison and a robust estimate of undercount.	Lengthy process to carry out and resource work.  Value added relies on the match-rate and ability to match on females.

In this work, we have demonstrated that information contained in free-text case notes can be used to supplement structured data for statistical aggregation. However, there is much more information contained in the text datasets which are yet to be explored. We have also shown the advantages of linking data across government departments (through a pilot matching exercise with HMRC) demonstrating the value of the BOLD programme.

## **6. Users**

The contents of this report will be of interest to government policy makers, the agencies responsible for offender management at both national and local levels,

providers, practitioners and others who want to understand more about children with a parent in prison and prisoners with children.

Government policymakers may also use these statistics to inform key elements of government policies. Offender management agencies may use these statistics to gain a better understand of the extent of the estimated number of children impacted by parental imprisonment. Key agencies include: HMPPS, private and voluntary sector providers of prison and probation services and local authorities.

## **7. Contact details and feedback on consultation**

You can send enquiries and feedback on these statistics to MoJ at [RR-pilot-BOLD@justice.gov.uk](mailto:RR-pilot-BOLD@justice.gov.uk)

For more information about the free-text modelling part of the project please contact [AI\\_for\\_linked\\_data@justice.gov.uk](mailto:AI_for_linked_data@justice.gov.uk).

*The free-text modelling work was supported by [The Alan Turing Institute](#) through a Turing Internship Network (TIN) internship.*