# Trust in data, and data in trust

Jim Knight, Labour Peer and Timo Hannay, SchoolDash

**Jim Knight, The Rt Hon. the Lord Knight of Weymouth**, works in education, digital technology and as a legislator. He is a director of Suklaa Ltd, providing advice to clients in education. Jim is a founder of xRapid, an AI diagnostic business. He is the Chair of E-Act Multi Academy Trust, the Digital Poverty Alliance and CAST. He is a board member of Century-Tech, MACAT International and GoBubble and sits on the advisory bodies for Nord Anglia, and BETT. As a government minister and MP, Jim's portfolios included rural affairs, schools, digital and employment. He was a member of Gordon Brown's Cabinet, before joining the Lords in 2010.

**Timo Hannay** is the founder of SchoolDash, an education data analytics firm that works with a wide range of collaborators, including media, publishers, edtech, charities, trusts and government. He is also a non-executive director of SAGE Publishing and of Arden University, and an advisor to Ada Lovelace Day and Maths4Girls.  He was previously the founding managing director of Digital Science and before that variously ran the online business of Nature Publishing Group, worked as a consultant at McKinsey & Company, wrote for The Economist, and investigated brain physiology at the University of Oxford and Waseda University in Tokyo.

Information is the currency of our age. Like more traditional forms of value, it can be used either to alleviate or exacerbate social ills. We believe that greater sharing and analysis of data is essential to more fully understand and address shortcomings in education. This must be done responsibly, benefiting the education system as a whole. One promising approach towards achieving this aim is the emerging concept of data trusts: legal entities that provide independent stewardship of data. This essay explores their potential in the context of education, particularly in the wake of the COVID-19 pandemic.

**The risks of data processing and usage**
In the online world, it is both the best and the worst of times. The internet has enabled access to information, opportunity and human connection in a way that was previously inconceivable.

Data is being generated, harvested and analysed at a scale that is transforming our economies and societies. To social scientists and policymakers, data provides a uniquely powerful observational tool – akin to the telescope for astronomers or the microscope for biologists. We no longer need to interview

1000 people and extrapolate; we can analyse the actual behaviour of whole populations in real time. In education, too, the adoption of online and artificial intelligence (AI)-driven approaches to learning has surged, especially in the wake of the COVID-19 pandemic. Indeed, we were perhaps lucky that it struck at a time when such alternatives to classroom-based teaching were even possible.

Yet any excitement at these remarkable and genuinely impressive developments must be tempered by deep concern about all sorts of adverse consequences.

Two years ago, the World Health Organization (WHO) declared that, alongside the COVID-19 pandemic, it was also fighting an 'infodemic' of misinformation (WHO, 2020). As we write, the ongoing war in Ukraine is being fought not only on land with guns and tanks, but also in cyberspace, with novel forms of propaganda and military intelligence. More than ever before, social media platforms have become geopolitical players and are struggling to tame the viral, runaway nature of their own networks and algorithms (Bushwick, 2022). It seems that every week brings news of yet another data leak or ransomware attack (Page, 2022). In the UK, the 5% of people who remain offline cite worries about privacy, identity theft and misuse of personal data as the most common reasons for foregoing the benefits of these new technologies (Lloyds Bank, 2021). The rise in online harms is prompting governments to seek to regulate large swathes of the internet through mechanisms such as the Online Safety Bill.

## Data in education

Like many other domains, education generates vast quantities of data, most of it held by private organisations, including educational technology (EdTech) firms, test publishers, tuition providers and survey companies. So far, these proprietary datasets have been mostly invisible and largely unexamined, even to the organisations that hold them.

The pandemic has started to change these attitudes to education data.

Datasets have grown even bigger, driven by the move to online learning. Their usefulness has become much more apparent, due to an urgent need to understand the effects of

lockdown and school closures on children's learning and wellbeing.

For example, the Education Policy Institute (EPI), an independent charity, was commissioned by the Department for Education (DfE) to analyse children's academic progress (or lack thereof) using data from Renaissance Learning, a commercial test provider (EPI, 2021). As schools closed during lockdown and then reopened, many of them used these tests to assess their pupils' individual progress and learning needs. By aggregating data from all such schools, the EPI was able to publish a series of studies during the course of the pandemic that characterised and quantified the national picture with respect to pupils' lost learning relative to previous cohorts. This stimulated the national debate and informed administrative decisions at every level, from central government to individual schools. It helped to answer such important questions as: exactly how much are pupils underperforming? How does this vary by age or geography? Which subjects and topics have been most affected?

On their own, such analyses do not solve problems, or even guarantee consensus about how to do so (indeed, there was considerable disagreement on this between the UK government and its own education recovery commissioner, leading to his eventual resignation; see Coughlan & Sellgren, 2021). But they were widely reported and discussed, and served the vital role of grounding the debate in objective reality rather than anecdote and preconception.

SchoolDash, a data analytics firm founded by one of us (TH), has been conducting similar work with another test publisher, Hodder Education (Hannay, 2021a). As well as providing alternative perspectives on the problem (since each dataset lent itself to slightly different analyses), a broad consistency between the EPI's results and those of SchoolDash helped to inspire confidence that the lost learning was real. In addition, SchoolDash has analysed data from EduKit (Hannay, 2020), a wellbeing survey company used by schools, and a number of online learning providers, including Oak National Academy (Hannay, 2021b), a government-funded initiative established during the pandemic to provide free online video lessons. These provided insights into how young people were

coping with home-based learning, both psychologically, in terms of social and emotional contentment, and practically, in terms of being able to access online lessons and engage effectively with the material provided.

These initiatives were largely products of COVID-19 and the associated lockdowns, when timely information about children's activities became especially important and official statistics were either too slow or, in some cases, absent altogether. But the potential of such analyses to help us understand and improve education goes well beyond the unique circumstances of the pandemic. If we can use such data to understand attainment gaps, technology divides and wellbeing deficits during lockdown, then why not during more normal times too? Educational inequalities and shortcomings are perennial challenges, and the mission to reduce them continues to deserve all the insight we can muster.

### The risks in education data uses vs. public trust

These developments are exciting – but potentially scary too. We know from the activities of big tech companies and totalitarian regimes that unfettered use of personal information can have bad effects, whether intended or not. Perhaps understandably, trust in governments and technology companies to use information responsibly is in decline in the UK (Wisniewski, 2020), the USA (PAC, 2021) and elsewhere. How can we enjoy the collective benefits while minimising the risk that important data might be used to serve narrow commercial or political interests rather than the interests of learners and the common good? This is particularly concerning when the data refers to individual people, and all the more so when those people are children.

The EU's General Data Protection Regulation (GDPR) has created a regime that reduces potential harm while maintaining many of the benefits that come from gathering and analysing personal data. However, it has not resulted in an increase in public trust (Wisniewski, 2020). Furthermore, insofar as this is based on user consent, it is difficult for parents or children (or anyone) to give consent for use that may be highly technical or somewhat uncertain in terms of outcome, as is often the case in analyses of subjects such

as educational disparities.

The Age Appropriate Design Code (AADC) further protects children's online activity in the UK, and this will soon be followed by an Online Safety Bill. These will help private and public sector technology providers to build on well-defined minimum standards. But additional, more flexible solutions are required to truly minimise risks and support the positive use of potentially sensitive information. It is abundantly clear that self-regulation by technology companies is not enough on its own. What other approaches might help?

Many kinds of organisation already use independent oversight to protect wider interests: schools have governors, charities have trustees and companies have non-executive directors. As data acquires increasing personal and societal relevance, perhaps it deserves similar safeguards too. This is the central idea behind 'data trusts', a relatively new concept that builds on existing data rights and trust law to provide independent fiduciary stewardship of data (ODI, 2018).

### The potential of data trusts in building trust in data

In October 2017, the UK Government published an independent review, *Growing the artificial intelligence industry in the UK* (Hall & Pesenti, 2017). It called for the 'development of data trusts, to improve trust and ease around sharing data'. Since then, interest and activity around this idea has steadily increased, with organisations such as the Open Data Institute (2018), the Alan Turing Institute (2019), the Ada Lovelace Institute (2021), Nesta (Mulgan & Straub, 2019) and the Data Trusts Initiative[1] (a collaboration between the Universities of Cambridge and Birmingham), all exploring and promoting the concept. They are now being talked about in commercial contexts, for example to represent users of a particular service, and their application in sensitive domains like AI (Mehonic, 2018) and healthcare (Milne et al., 2020) is being actively pursued.

Data trusts can take a wide variety of forms (O'Hara, 2019). A common model is for the trustees to represent the interests of a well-defined group of people, such as a local community or a cohort that is the subject of a research project[2] (somewhat akin to an academic ethical review board). We believe that they

can also be useful in representing the wider interests of, say, the education system as a whole, and even society at large.

With this in mind, we are collaborating with other individuals and organisations to establish an education data trust[3] that will aggregate information across multiple proprietary data holders and public sources, using these to conduct analyses that will provide actionable insights for participants at all levels of the education system. This builds on the work conducted before and during the pandemic described above, and will help to answer questions such as: how well have children caught up on learning following the reopening of schools? What are the current demographic and socioeconomic disparities in attainment? How does access to technology vary by pupil age or location? How big a role is online education playing, and how is this evolving over time? These are just a few examples from an almost endless list of possibilities. Sounder, more timely answers to these questions would support better-informed public debate and more effective education policies, ultimately to the benefit of children.

Crucially, all of these activities will be overseen by a board of independent, knowledgeable and highly regarded trustees, who will be tasked with representing the interests of the data subjects and the wider education system. Unless they approve of a particular activity, it will not be allowed to happen. This will obviously constrain our ability to conduct whatever analyses we (and the data providers) choose to conduct. But, far from holding us back, we expect this to be an enabler. Done properly, a rigorous and transparent approach to data oversight will instil greater confidence in our activities, leading to more, not fewer, opportunities to access new sources of information and use them to derive valuable insights. It will initially focus on schools in England and other parts of the UK, but with the potential to extend into further or higher education, and into other territories. We also hope that it will serve as a template for those who would like to establish similar initiatives in other domains, just as we have been inspired by emerging projects elsewhere (see, for example, the Data Trust Initiative's 2022 pilot projects[4]).

Data trusts are not a substitute for legislation, public regulation or even industry self-regulation, but they are an important addition to the mix. Neither are they a fool-proof, catch-all solution, just as the existence of a governing body does not guarantee the proper running of a school. But oversight by a group of independent, knowledgeable and credible trustees who are answerable not to senior executives or shareholders but to data subjects and wider society surely represents a step in the right direction. Indeed, not to make such a move towards better stewardship and greater transparency is to invite further scepticism towards data analysis of any kind, decreasing its potential in solving real-world problems and letting go a huge opportunity to pursue the common good.

We expect that there will be further applications of data trusts that go beyond the kinds of analyses described here. The use of AI in education is a particularly pressing area. It offers huge opportunities by helping teachers to better match pedagogy and resources to individual learners, and to do so far more efficiently than has ever been possible before. But it requires vast quantities of training data, which, in turn, raises ethical questions.

An instructive example is the use by Google's DeepMind of the health records of 1.5 million patients at an NHS Trust (BBC News, 2021). The company used these to train an AI system that detects people at risk of kidney injury. Beyond the concerns of the Information Commissioner around data privacy, there was also controversy over the intellectual and commercial rights, since NHS patients had shared data with a public service provider only to see it being used to create a commercial product that was then sold back to the NHS at taxpayers' expense.

Similar conflicts of interest could arise in schools. For example, if a company uses one product to collect learner data and then applies it in developing a second product, how should the resulting value be shared between the company and the school, and who is looking out for the interests of individual learners? There are no easy or universal answers to such questions, which makes it all the more important to establish systems of oversight that take adequate account of the full range of interests involved.

## Towards responsible data uses

Despite increased regulation, trust in technology is in decline and this backlash will continue to constrain the adoption of new, potentially beneficial innovations. While not a panacea, data trusts represent an important part of the solution. We anticipate that they will become a standard means to ensure responsible use of data across education and beyond, especially where the information concerns children or other potentially vulnerable groups. What Creative Commons has done for content sharing, and Wikipedia has done for knowledge dissemination, data trusts might yet do for online information. By harnessing the power of openness and collaboration, we hope that they will help to support the internet as a force for individual and collective good.

Ada Lovelace Institute & AI Council. (2021). Exploring legal mechanisms for data stewardship, 4 March

Alan Turing Institute, The (2019). Data trusts workshop

BBC News. (2021). DeepMind faces legal action over NHS data use, 1 October

Bushwick, S. (2022). Russia's information war is being waged on social media platforms. *Scientific American*, 8 March

Coughlan, S., & Sellgren, K. (2021). School catch-up tsar resigns over lack of funding. BBC News, 2 June

EPI (Education Policy Institute). (2021). *EPI research for the Department for Education on pupil learning loss*. 29 October

Hall, W., & Pesenti, J. (2017). *Growing the artificial intelligence industry in the UK*. Department for Digital, Culture, Media & Sport and Department for Business, Energy & Industrial Strategy

Hannay, T. (2020). The lockdown experiences of pupils in England. SchoolDash, 27 August

Hannay, T. (2021a). The effects of educational disruption on primary school attainment in summer 2021. SchoolDash, 31 August

Hannay, T. (2021b). What Oak National Academy usage tells us about education during the pandemic. SchoolDash, 5 November

Lloyds Bank. (2021). *UK Consumer Digital Index 2021*

Mehonic, A. (2018). Can data trusts be the backbone of our future AI ecosystem? The Alan Turing Institute, 3 October

Milne, R., Sorbie, A., & Dixon-Woods, M. (2020). What can data trusts for health research learn from participatory governance in biobanks? *Journal of Medical Ethics, 48*(5)

Mulgan, G., & Straub, V. (2019). The new ecosystem of trust. Nesta, 21 February

ODI (Open Data Institute). (2018). What is a data trust? 10 July

O'Hara, K. (2019). *Data trusts: Ethics, architecture and governance for trustworthy data stewardship*. WSI White Paper # 1, February. Web Science Institute

PAC (Public Affairs Council). (2021). *2021 Public Affairs Pulse Survey report*

Page, C. (2022). Thousands of Nvidia employee passwords leak online as hackers' ransom deadline looms. TechCrunch, 4 March

WHO (World Health Organization). (2020). Munich Security Conference Speech. WHO Director-General, 15 February

Wisniewski, G. (2020). Losing faith: The UK's faltering trust in tech. Edelman, 30 January

1   See: https://datatrusts.uk

2   For more information, visit: https://datatrusts.uk/pilot-projects-1

3   More information about a project concerning education data trust will soon be available at: https://educationdatatrust.com

4   For more information, visit: https://datatrusts.uk/pilot-projects-1