


**Please cite the Published Version**

Aykol, Muratahan, Hummelshøj, Jens S, Anapolsky, Abraham, Aoyagi, Koutarou, Bazant, Martin Z, Bligaard, Thomas, Braatz, Richard D, Broderick, Scott, Cogswell, Daniel, Dagdelen, John, Drisdell, Walter, Garcia, Edwin, Garikipati, Krishna, Gavini, Vikram, Gent, William E, Giordano, Livia, Gomes, Carla P, Gomez-Bombarelli, Rafael, Balaji Gopal, Chirranjeevi, Gregoire, John M, Grossman, Jeffrey C, Herring, Patrick, Hung, Linda, Jaramillo, Thomas F, King, Laurie , Kwon, Ha-Kyung, Maekawa, Ryosuke, Minor, Andrew M, Montoya, Joseph H, Mueller, Tim, Ophus, Colin, Rajan, Krishna, Ramprasad, Rampi, Rohr, Brian, Schweigert, Daniel, Shao-Horn, Yang, Suga, Yoshinori, Suram, Santosh K, Viswanathan, Venkatasubramanian, Whitacre, Jay F, Willard, Adam P, Wodo, Olga, Wolverton, Chris and Storey, Brian D (2019) The Materials Research Platform: Defining the Requirements from User Stories. *Matter*, 1 (6). pp. 1433-1438. ISSN 2590-2385

**DOI:** <https://doi.org/10.1016/j.matt.2019.10.024>

**Publisher:** Elsevier BV

**Version:** Accepted Version

**Downloaded from:** <https://e-space.mmu.ac.uk/624624/>

**Usage rights:**  [Creative Commons: Attribution-Noncommercial-No Derivative Works 4.0](#)

**Additional Information:** This is an Author Accepted Manuscript of an article published in *Matter* by Elsevier.

**Enquiries:**

If you have questions about this document, contact [openresearch@mmu.ac.uk](mailto:openresearch@mmu.ac.uk). Please include the URL of the record in e-space. If you believe that your, or a third party's rights have been compromised through this document please see our Take Down policy (available from <https://www.mmu.ac.uk/library/using-the-library/policies-and-guidelines>)

## The Materials Research Platform: Defining the Requirements from User Stories

Muratahan Aykol<sup>1\*</sup>, Jens Hummelshøj<sup>1</sup>, Abraham Anapolsky<sup>1</sup>, Koutarou Aoyagi<sup>2</sup>, Ryoji Asahi<sup>3</sup>, Martin Bazant<sup>4</sup>, Thomas Bligaard<sup>5</sup>, Richard Braatz<sup>4</sup>, Scott Broderick<sup>6</sup>, Daniel Cogswell<sup>4</sup>, John Dagdelen<sup>7</sup>, Walter Drisdell<sup>7</sup>, Edwin Garcia<sup>8</sup>, Krishna Garikipati<sup>9</sup>, Vikram Gavini<sup>9</sup>, William Gent<sup>10</sup>, Livia Giordano<sup>4</sup>, Carla Gomes<sup>11</sup>, Rafael Gomez-Bombarelli<sup>4</sup>, Chirranjeevi Gopal<sup>1</sup>, John Gregoire<sup>12</sup>, Jeffrey Grossman<sup>4</sup>, Patrick Herring<sup>1</sup>, Linda Hung<sup>1</sup>, Thomas Jaramillo<sup>10</sup>, Laurie King<sup>10</sup>, Ha-Kyung Kwon<sup>1</sup>, Ryosuke Maekawa<sup>2</sup>, Andrew Minor<sup>13</sup>, Joseph Montoya<sup>1</sup>, Tim Mueller<sup>14</sup>, Colin Ophus<sup>7</sup>, Krishna Rajan<sup>6</sup>, Rampi Ramprasad<sup>15</sup>, Brian Rohr<sup>10</sup>, Daniel Schweigert<sup>1</sup>, Yang Shao-Horn<sup>4</sup>, Yoshinori Suga<sup>2</sup>, Santosh Suram<sup>1</sup>, Venkat Viswanathan<sup>16</sup>, Jay Whitaker<sup>16</sup>, Adam Willard<sup>4</sup>, Olga Wodo<sup>6</sup>, Chris Wolverton<sup>17</sup>, Brian Storey<sup>1\*</sup>

<sup>1</sup>Toyota Research Institute, United States

<sup>2</sup>Toyota Motor Corporation, Japan

<sup>3</sup>Toyota Central R&D Labs., Inc., Japan

<sup>4</sup>Massachusetts Institute of Technology, United States

<sup>5</sup>SLAC National Accelerator Laboratory, United States

<sup>6</sup>University at Buffalo, United States

<sup>7</sup>Lawrence Berkeley National Laboratory, United States

<sup>8</sup>Purdue University, United States

<sup>9</sup>University of Michigan, United States

<sup>10</sup>Stanford University, United States

<sup>11</sup>Cornell University, United States

<sup>12</sup>California Institute of Technology, United States

<sup>13</sup>University of California Berkeley, United States

<sup>14</sup>Johns Hopkins University, United States

<sup>15</sup>Georgia Tech, United States

<sup>16</sup>Carnegie Mellon University, United States

<sup>17</sup>Northwestern University, United States

\*Corresponding Authors: Muratahan Aykol ([murat.aykol@tri.global](mailto:murat.aykol@tri.global)) or Brian Storey ([brian.storey@tri.global](mailto:brian.storey@tri.global))

**Summary: A recent meeting focused on** accelerated materials design and discovery examined user requirements for a general, collaborative, integrative and on-demand *materials research platform*.

What common elements are necessary to create a general framework for materials innovation? Here, we provide a retrospective analysis of high-level themes that emerged from a focused discussion on the requirements for a future **materials research platform**. These discussions occurred during the annual Toyota Research Institute - Accelerated Materials Design and Discovery meeting (May 29, 2019, in Boston, MA) with nearly 50 field experts from universities, US national laboratories, Toyota Research Institute in the United States, and Toyota Motor Corporation in Japan. These researchers contributed ideas towards what capabilities such a platform should have to accelerate the design and discovery of materials.

To maximize the information captured in the form of these field expert's ideas, we followed a strategy comprised of exploitation and exploration inspired by knowledge acquisition strategies in machine-learning (ML). In a first session, we matched researchers with themes that we expect them to be knowledgeable about (hence exploitation). In a follow-up session, we did a quasi-random assignment of researchers to themes (hence exploration), with the goal of capturing unique ideas that might elude experts embedded deep in their work. Sixteen preselected themes relevant for a materials research platform<sup>a</sup> were discussed twice in two unique groups, one formed as an exploitation and the other formed as an exploration team. We digitally recorded the ideas in the form of "user stories" from agile software development practices. The meeting captured many stories that would be considered obvious, but several unexpected ideas related to the platform emerged. An initial parsing revealed several major interrelated themes for the design of a useful platform: (i) data and knowledge assets, (ii) automation of science and (iii) integrative approaches, as outlined in Figure 1.

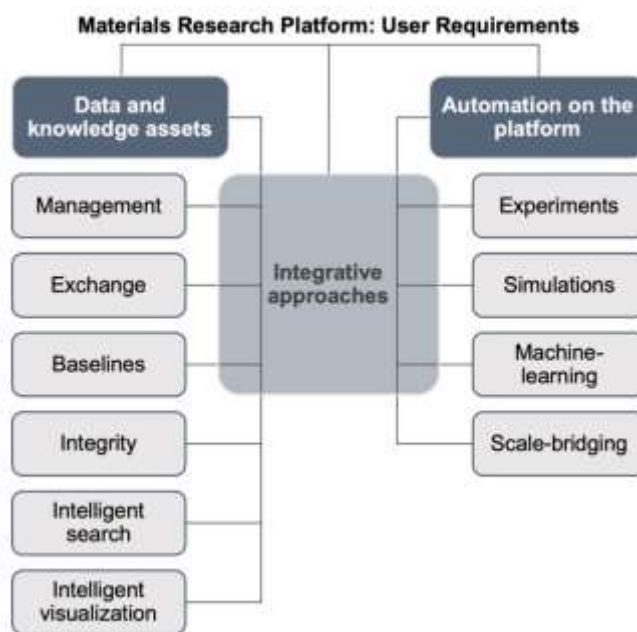
The ideas related to data and knowledge assets rely on the FAIR principles,<sup>1</sup> but with added distinguishing capabilities relevant for a materials research platform. For example, the platform should enable sharing and collaboration, not just around data, but also knowledge assets such as machine learning models or scientific workflows. Baselines to gauge new findings are critical, but often overlooked. Artificial intelligence (AI) assisted search and visualization could amplify the scientific abilities of human researchers. As part of the automation, experimental workflows are envisioned to be run *on-demand*, where tasks are picked up by relevant laboratories with fully or partially automated experimental capabilities, forming collaborative networks via the platform. Given the state-of-the-art for automation in modeling and simulation, a similar but more "productized" automated capability with web-based user interfaces is envisioned to assist researchers to run *on-demand* simulations complementary to experiments, or train or use, *on-demand* machine-learning models, without deep knowledge in any particular method. Automation of scale-bridging, which would entail at least designing workflows for codes beyond density functional theory (DFT) and multi-scale application programming interfaces (APIs) for linking the codes, emerged as a game-changing capability to bridge the gap between the computer-design or laboratory and device level properties, but is also deemed one of the major outstanding scientific challenges. The design of integrative approaches on the platform that

---

<sup>a</sup> There were 16 themes selected: Adaptive systems: active-learning & beyond, Automation of experiments, Automation of simulations, Collaboration, Data ingestion and sharing, Integration, Knowledge discovery, Machine-learning for experiments, Machine-learning for simulations, Multi-fidelity & uncertainty quantification, Reproducibility & provenance, Scale-bridging, Simulation tools, Software infrastructure, Text mining & NLP and Visualization.

leverage diverse datasets, physical, empirical, computational or statistical models, and experiments would require a modularized software framework, cost estimation functions, uncertainty quantification and fidelity assessments of new data points, and subsequently, discovery “engines” like optimal experiment design (OED), active-learning or optimization can be built for diverse material-device domains, from batteries to catalysis.

Materials science of the future is expected to be interwoven with data, automation, machine-learning and other emerging information technologies. Many aspects of these paradigms are being actively reviewed, debated and discussed by the materials community.<sup>2–4</sup> In this article, we focus on the requirements for the enablement of a useful, general, next-generation materials research platform that would combine, build and expand on these data-driven paradigms to enable discovery and innovation. A number of academic and industry teams are actively engaged in efforts relevant for this vision.<sup>5–11</sup> We expect that the entire materials community will benefit from the distilled summary of ideas and requirements of the envisioned future system that we present in this article. The following sections expand the concepts underlying the platform related themes of data and knowledge assets, automation, and integration for materials research.



**Figure 1.** Requirement categories identified by the researchers as the main pillars of a useful materials research platform: data and knowledge assets, automation on the platform and integrative approaches.

### The core of the platform: data and knowledge assets

A research platform inherits, generates, stores and serves data and knowledge as part of its mission. Such a platform’s utility for the scientist, therefore, is tied with how these key bits of information are assimilated and managed, how their quality, reliability and integrity are judged,

how their exchange and dissemination are enabled, and ultimately whether the informatics aspect of the platform is meeting the needs of the scientist and helping them innovate.

**Data management:** The requirements regarding data management are a subset of the now well-accepted FAIR principles: *Findable, Accessible, Interoperable, Reusable*.<sup>1</sup> As a basic requirement, the platform is expected to contain standardized datasets. Standardization of all data imported to, created on, and disseminated by the platform would entail adoption of established ingestion procedures, data formats, capturing of metadata (e.g. experimental conditions), provenance and instrument logs, for instance to differentiate human vs machine generated data. In addition, the platform is expected to deliver not only well-known computational databases, but also diverse, large, high-quality *experimental* datasets, where inclusion of “dark data” (i.e. data that is considered a “negative” result and not publishable) is essential. Inclusion of all data is key to removing human bias from datasets. Interaction with the data system needs to be easy and intuitive, programmatically or via a web-based user interface (UI) that allows easy or automated upload, instant visualization, search and sharing, and is connected to the other components of the platform. Importance of a *simple yet powerful UI* for all components of the platform constantly came up throughout the discussions: we will not repeat that requirement and it should be assumed by default to be a core component for every module hereafter.

**Collaboration and knowledge exchange:** Data sharing is a core component of today’s data-driven research. The storage and sharing of analytical tools, machine-learning models, workflows and other knowledge (overall, knowledge assets) as well as experimental resources (for instance, see Automation of Experiments) via the platform would enable a more complete, collaborative research experience. The system is envisioned to enable citations for all such shareables, and provide utilization logs, allowing the apportioning of separate credit to datasets, models, and scientific results.<sup>9</sup> Furthermore, the platform may allow rapid user feedback and community review, which will motivate developers to maintain quality components on the platform (e.g., clearly licensed, well-documented, and maintained code repositories). Given the diversity of the materials research community, the platform is expected to balance certain users’ and institutions’ desire for privacy while incentivizing the sharing of methods, data, and scientific results. As such, some interactions may be open collaborations or crowdsourcing, others might require the platform to act as a marketplace or exchange, while still others can have a training or educational component.

**Baselines:** Baselines help gauge where a new scientific finding stands, but are often lacking. The data system on the platform can host curated baseline materials, with all of their properties, history, cost (e.g. established Li-ion electrodes, electrolytes, established fuel-cell membranes, catalysts, etc.) such that new material discoveries can be compared to the current state of the art. Relevant design parameters for finding replacement materials can be accessed easily. The same is true for analytical methods and tools, which would require having standardized, benchmark datasets (e.g. curated DFT datasets) and baseline models trained on them.

**Data integrity tools:** Integrity of data is critical and can be enabled in part by standardization and ingestion procedures. More advanced systems, which are often not part of existing scientific workflows (experiments in particular), such as data validation and anomaly detection can reinforce data integrity. Such capabilities, if they operate on the fly, can increase the value of the data and can make the platform attractive for research groups that produce live data streams. Anomaly or outlier detection, as an unsupervised method, can also play a complementary role in discovery, as discoveries often appear as outliers.

**Intelligent search:** Reusability and discoverability of data on the platform are expected by default. A feature that emerged as part of multiple themes is an “Intelligent Search” system for materials research. The system is envisioned to operate beyond chemical formula or material labels, and can search over properties, models (e.g. machine-learning models), methods (e.g. synthesis recipes, characterization methods), tools and other metadata.<sup>12</sup> A search capability that is based on data itself (e.g. similarity match with a user supplied XPS spectrum) would be a game changer. More context aware, interactive search capability beyond simple keywords should be implemented (e.g. “Find alternatives to sol-gel synthesis of  $\text{Cr}_2\text{O}_3$ ”). In general, the system could recommend new research directions including; materials, techniques, keywords, collaborators, or publications. Specialized recommender systems, as described above for materials, keywords, collaborators, publications, or simulation tools came up in distinct groups.

**Intelligent visualization:** A powerful, scientifically focused visualization technology can be considered a distinguishing feature of a next-generation platform. A useful capability would include standard, automated exploratory data analysis and visualization for data on the platform. Visualizations should extend beyond current workflows (e.g. plotting experimental or computational data, and data derived from those) to visualization of high-dimensional parameter spaces (e.g. embeddings in machine-learning models), visualizing relationships between codes, models and simulations (e.g. see scale-bridging) or mapping data provenance and property relationships, where graph/network based representations may be useful.

## Automation on the platform

Generation of high-quality, large-volume and consistent data streams is often enabled by automation of manual tasks, making processes less prone to human error and increasing throughput. Automation of experiments and simulations are two fundamental paradigms that were considered, where both converged to a desired “on-demand” capability on the platform. Automation of other components, such as machine-learning and scale-bridging independently emerged.

**Automation of experiments:** Automation of experiments is expected to provide critical functionalities for materials discovery such as reduction of human bias and enabling rapid access to multiple material design axes such as composition, reproducibility, or processing. As a basic functionality, the platform is expected to be integrated with a distributed system of experimental facilities, to connect to their data streams and to enable experiment requests.<sup>10</sup> The platform should provide capabilities for creation, execution and moderation of workflows

that span one or more experimental facilities, and also should recommend such workflows for specific applications, experimental cost estimates, and fidelity. The platform can provide *on-demand experiments*, a marketplace for experiments or a collaborative closed-network, where a user can for example, request a synthesis of a target material or characterization of a sample at the participating facilities or labs. Importance of automation of low-throughput, repetitive experiments with the aid of robotics was also highlighted.

**Automation of simulations:** Computer simulations are, by their nature, more amenable to automation than experiments. Automation of DFT formed the seed for the present era of data-driven materials science by providing large, reliable material datasets.<sup>5-7</sup> Hence, the stories related to automation of simulation focused on capabilities beyond automation of DFT itself: such as molecular dynamics (MD), coarse-graining methods, phase field, and beyond, to predict macroscopic and device level properties. In addition, in analogy with “on-demand experiments”, a paradigm of “on-demand simulations” emerged, where the platform can provide an easy to use interface for users (simple enough to be useful to non-specialists) to request new simulations complementary to their ongoing experiments. A recommender system for types of simulations, parameters, and ready-to-use license arrangements would add value. As mentioned before, the platform should display relevant benchmarks for simulation tools (e.g. accuracy, performance, cost) and document use cases.

**Automation of machine-learning:** An easy to use machine-learning and analytics module on the platform, backed by a powerful UI that requires no deep expertise, developed as a common feature desired among multiple discussion groups. To create the necessary input for training predictive models, automated featurization of materials (or other entities) should be a part of this module. In addition, for more advanced practitioners and more complex predictive problems, a comprehensive machine-learning arsenal can be provided: e.g. for image processing, spectroscopy, natural language processing (NLP), deep learning, machine-learning for rare-events, failures, stochastic events, material-processing relationships, microstructure-property relationships, physically-informed ML models, noisy data, as well as generative and evolutionary models. The UI should display convenient visual information, such as performance metrics, feature importance in models, and system should alert the user when there are concerns about the data integrity or quality (bias, anomalies, etc. See “Data integrity tools”). In addition, unsupervised methods that identify relationships in the data and/or capture low dimensional representations should be available.

**Automation of scale-bridging:** Scale-bridging is required for a more complete assessment of device-level properties of material systems. Often a small change in material properties used as part of a device, requires redesign or re-evaluation of many other components of the same device. Incorporation of new materials in devices has traditionally been a long, slow and costly process. Today scale-bridging is still a major roadblock in materials research, and mostly performed on an *ad hoc* basis. The need for automated scale bridging was strongly emphasized and also acknowledged as a major scientific challenge. As the most basic functionality, a visual relationship between simulation techniques that can operate at multiple scales has the potential to guide the users towards a hand-crafted scale-bridging study. It was noted that the data

transfer between different scales is not sequential or one-directional, and there can often be data transferred from all scales to the others (e.g. DFT to device level, DFT to MD, MD to device level, DFT to finite element, and so on.) and the transfer can be bi-directional (e.g. device design informs phase-field or phase field informs device design).

For effective scale-bridging, one should parameterize and automate simulations beyond DFT and describe the “contracts” and dependencies between inputs and outputs of such simulations. Such contracts can be framed as “Scale-bridging APIs” where input property requirements for methods are documented and codified, to enable programmatic integration between simulation software that operates at different length and time scales. As mentioned above, a graph/network of data transfers and dependencies of simulation tools can be constructed.

### **Integrative approaches on the platform**

The research envisioned to be enabled on the platform requires blending or integration of many components, tools and/or datasets. Several such paradigms emerged from the user stories that are centered on integration, where the ability to be automated and modularized was an expected quality for all relevant tools. Data fusion, where multiple datasets are combined to enrich the existing data, and to create new datasets, is the most basic example. Scale-bridging is a fundamental integration challenge highlighted in previous section.

An emerging paradigm for discovery of materials and processes is the application of cyclic, active-learning, optimization or OED based feedback-loop systems, where the science (and the underlying decision making) itself is partially automated.<sup>4</sup> The materials research platform should provide modular, plug-and-play, automatable closed-loop capability to enable this form of research. This capability needs to be easy to integrate with both experimental and computational data streams. In addition to the software and analytical infrastructure required for the process, cost estimating functions for experimental and simulation processes emerged as a key feature to have on the platform.

Uncertainty quantification is essential for automated, closed-loop research systems. Experts highlighted the importance of having uncertainty estimates or confidence intervals available on all experimental and computational measurements, parameters, and outputs. How uncertainty propagates as data is being transformed (e.g. in machine-learning or scale-bridging), and how that affects the reliability of resulting predictions also remains an open question. Experts mentioned potential benefit from incorporation of information-theoretic approaches into the system (e.g. information gain) as well as availability of classification tables for fidelities and costs of acquiring experimental or computational data points or those from surrogate models (machine-learning models, empirical models, etc.). Ultimately, multi-fidelity optimization where uncertainties, fidelities and cost are taken into account offer a viable, general pathway for integration of computational and experimental data generation pipelines to solve complex scientific problems. These features should have a presence on the platform as modular systems.



## Citations

- <sup>1</sup> M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.W. Boiten, L.B. da Silva Santos, P.E. Bourne, J. Bouwman, A.J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C.T. Evelo, R. Finkers, A. Gonzalez-Beltran, A.J.G. Gray, P. Groth, C. Goble, J.S. Grethe, J. Heringa, P.A.C. t Hoen, R. Hooft, T. Kuhn, R. Kok, J. Kok, S.J. Lusher, M.E. Martone, A. Mons, A.L. Packer, B. Persson, P. Rocca-Serra, M. Roos, R. van Schaik, S.A. Sansone, E. Schultes, T. Sengstag, T. Slater, G. Strawn, M.A. Swertz, M. Thompson, J. Van Der Lei, E. Van Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg, K. Wolstencroft, J. Zhao, and B. Mons, *Sci. Data* **3**, 1 (2016).
- <sup>2</sup> K. Alberi, M.B. Nardelli, A. Zakutayev, L. Mitas, S. Curtarolo, A. Jain, M. Fornari, N. Marzari, I. Takeuchi, M.L. Green, M. Kanatzidis, M.F. Toney, S. Butenko, B. Meredig, S. Lany, U. Kattner, A. Davydov, E.S. Toberer, V. Stevanovic, A. Walsh, N.-G. Park, A. Aspuru-Guzik, D.P. Tabor, J. Nelson, J. Murphy, A. Setlur, J. Gregoire, H. Li, R. Xiao, A. Ludwig, L.W. Martin, A.M. Rappe, S.-H. Wei, and J. Perkins, *J. Phys. D. Appl. Phys.* **52**, 013001 (2019).
- <sup>3</sup> D.P. Tabor, L.M. Roch, S.K. Saikin, C. Kreisbeck, D. Sheberla, J.H. Montoya, S. Dwaraknath, M. Aykol, C. Ortiz, H. Tribukait, C. Amador-Bedolla, C.J. Brabec, B. Maruyama, K.A. Persson, and A. Aspuru-Guzik, *Nat. Rev. Mater.* **3**, 5 (2018).
- <sup>4</sup> P. V Balachandran, D. Xue, J. Theiler, J. Hogden, and T. Lookman, *Sci. Rep.* **6**, 19660 (2016).
- <sup>5</sup> A. Jain, S.P. Ong, G. Hautier, W. Chen, W.D. Richards, S. Dacek, S. Cholia, D. Gunter, D. Skinner, G. Ceder, and K.A. Persson, *APL Mater.* **1**, 011002 (2013).
- <sup>6</sup> S. Kirklin, J.E. Saal, B. Meredig, A. Thompson, J.W. Doak, M. Aykol, S. Rühl, and C. Wolverton, *Npj Comput. Mater.* **1**, 15010 (2015).
- <sup>7</sup> S. Curtarolo, W. Setyawan, S. Wang, J. Xue, K. Yang, R.H. Taylor, L.J. Nelson, G.L.W. Hart, S. Sanvito, M. Buongiorno-Nardelli, N. Mingo, and O. Levy, *Comput. Mater. Sci.* **58**, 227 (2012).
- <sup>8</sup> C. Draxl and M. Scheffler, *J. Phys. Mater.* **2**, 036001 (2019).
- <sup>9</sup> B. Blaiszik, K. Chard, J. Pruyne, R. Ananthakrishnan, S. Tuecke, and I. Foster, *JOM* **68**, 2045 (2016).
- <sup>10</sup> M.L. Green, C.L. Choi, J.R. Hattrick-Simpers, A.M. Joshi, I. Takeuchi, S.C. Barron, E. Campo, T. Chiang, S. Empedocles, J.M. Gregoire, A.G. Kusne, J. Martin, A. Mehta, K. Persson, Z. Trautt, J. Van Duren, and A. Zakutayev, *Appl. Phys. Rev.* **4**, (2017).
- <sup>11</sup> <https://citrite.io>, <https://www.schrodinger.com>, <https://grantadesign.com>, <https://exabyte.io>, and <https://www.kebotix.com>, (n.d.).
- <sup>12</sup> V. Tshitoyan, J. Dagdelen, L. Weston, A. Dunn, Z. Rong, O. Kononova, K.A. Persson, G. Ceder, and A. Jain, *Nature* **571**, 95 (2019).