



Main Manuscript for

Artificial Intelligence in Food and Nutrition Evidence: The Challenges and Opportunities based on a Convening of Content Experts

^{1,2}Regan L. Bailey, ^{1,3}Amanda J. MacFarlane, ⁴Martha S. Field, ^{5,6}Ilias Tagkopoulos, ⁷Sergio E. Baranzini, ⁸Kristen M. Edwards, ⁹Christopher J. Rose, ¹⁰Nicholas J. Schork, ¹¹Akshat Singhal, ¹²Byron C. Wallace, ²Kelly P. Fisher, ⁵Konstantinos Markakis, ¹Patrick J. Stover

¹Department of Nutrition, Texas A&M University, College Station Texas, 77845 USA.

²Texas A&M University, Institute for Advancing Health Through Agriculture, College Station Texas, 77845 USA.

³Texas A&M Agriculture, Food, and Nutrition Evidence Center, Fort Worth, Texas, 76102 USA.

⁴Cornell University, Division of Nutritional Sciences, Ithaca NY 14850 USA.

⁵Department of Computer Science and Genome Center, University of California, Davis, Davis, California 95616 USA

⁶USDA/NSF AI Institute for Next Generation Food Systems (AIFS), University of California, Davis, Davis, California 95616 USA.

⁷Department of Neurology, Weill Institute for Neurosciences, University of California, San Francisco, San Francisco, CA 94158 USA.

⁸Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139 USA.

⁹Cluster for Reviews and Health Technology Assessments, Norwegian Institute of Public Health, Oslo, Norway and Centre for Epidemic Interventions Research (CEIR), Norwegian Institute of Public Health, Oslo, Norway.

¹⁰Translational Genomics Research Institute (TGen), City of Hope National Medical Center, Phoenix, AZ, 85027 USA.

¹¹Department of Computer Science and Engineering, University of California, San Diego, La Jolla, CA, 92093 USA.

¹²Khoury College of Computer Sciences, Northeastern University, Boston, Massachusetts, USA.

¹³Texas A&M University, 801 Cherry Street, Suite 850, Fort Worth, Texas 76102 USA.

Corresponding Author: Patrick J. Stover

Patrick.stover@tamu.edu

Texas A&M University

801 Cherry Street, Suite 850

Fort Worth, Texas, 76102

© The Author(s) 2024. Published by Oxford University Press on behalf of National Academy of Sciences. This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact reprints@oup.com for reprints and translation rights for reprints. All other permissions can be obtained through our RightsLink service via the Permissions link on the article page on our site—for further information please contact journals.permissions@oup.com.

1
2 **Disclaimer:** The findings and conclusions in this report are those of the authors and should not
3 be construed to represent any official National Academies of Sciences, Engineering, and
4 Medicine, USDA, HHS, or U.S. Government determination or policy nor necessarily represent the
5 view of meeting participants.
6

7 **Author Contributions:** RLB and PJS drafted the manuscript informed by expert meeting
8 speakers AJM, MSF, IT, SEB, KME, CJR, NJS, AS, and BCW. Expert meeting speakers revised
9 the draft manuscript and provided figures. KPF provided content and managed all revisions. KM
10 provided technical advice.

11 **Competing Interest Statement:** SEB is co-founder of Mate Bioservices; I.T. is the founder of
12 PIPA LLC, an AI company on food, nutrition, and health.

13 **Classification:** Biological Sciences, Physical Sciences.

14 **Keywords:** nutrition, computed evidence, artificial intelligence, evidence synthesis, systematic
15 reviews.

16 **This PDF file includes:**

17 Main Text
18 Figures 1 to 4
19

20 **Abstract**

21 Science-informed decisions are best guided by the objective synthesis of the totality of evidence
22 around a particular question and assessing its trustworthiness through systematic processes.
23 However, there are major barriers and challenges that limit science-informed food and nutrition
24 policy, practice, and guidance. First, insufficient evidence, primarily due to acquisition cost of
25 generating high-quality data, and the complexity of the diet-disease relationship. Furthermore, the
26 sheer number of systematic reviews needed across the entire agriculture and food value chain,
27 and the cost and time required to conduct them, can delay the translation of science to policy.
28 Artificial intelligence (AI) offers the opportunity: 1) to better understand the complex etiology of
29 diet-related chronic diseases; 2) to bring more precision to our understanding of the variation
30 among individuals in the diet-chronic disease relationship; 3) to provide new types of computed
31 data related to the efficacy and effectiveness of nutrition/food interventions in health promotion;
32 and 4) to automate the generation of systematic reviews that support timely decisions. These
33 advances include the acquisition and synthesis of heterogeneous and multimodal datasets. This
34 perspective summarizes findings from a meeting of experts convened at the National Academy of
35 Sciences, Engineering and Medicine. The purpose of the meeting was to examine the current
36 state and future potential of AI in generating new types of computed data as well as automating
37 the generation of systematic reviews to support evidence-based food and nutrition policy, practice
38 and guidance.

1 Main Text

2 Introduction

3 Science-informed decisions are guided by the objective synthesis of the totality of evidence
4 around a particular question, and assessing its trustworthiness through the process of conducting
5 a systematic review. This approach has become fundamental to evidence-based food and
6 nutrition policy, practice, and guidance (1-3). Evidence synthesis and evaluation considers the
7 strength of all forms of scientific data and is used across medicine, public health, and the social
8 sciences.

9 Systematic reviews (SR) guide the process for setting essential nutrient intake recommendations
10 for individuals and populations, such as the Dietary Reference Intakes (4), and for food-based
11 intake recommendations including the Dietary Guidelines for Americans (3). Guidance on nutrient
12 and other food substances is based on derived normative values and include the Recommended
13 Dietary Allowance, the Estimated Average Requirement, and the Tolerable Upper Intake Level.
14 The DRIs inform food and nutrition policies including the Dietary Guidelines for Americans (5),
15 food fortification policies (6), food assistance programs (7), food safety, labelling and other
16 regulatory decisions (8), nutrition education programs (9), and can influence food production
17 systems (8). The food and agriculture economy contributes 1.53 trillion dollars to the United
18 States gross domestic product (~5.6 % of overall (10)), while food-related health effects due to
19 cardiometabolic diseases including hypertension, stroke, type 2 diabetes, and heart disease,
20 account for 50 billion USD/year in health care costs (11). highlighting the importance of bringing
21 the very best and current science available to policy and other decision makers. However, there
22 are major barriers and bottlenecks that limit the opportunity to achieve science-informed food and
23 nutrition policy. These include a dearth of high-quality scientific data to inform policy decisions,
24 the costs of generating high-quality food and nutrition experimental data, and the vast and rapidly
25 growing literature base; the sheer number of systematic reviews required to address all policy-
26 related questions across the entire agriculture and food value chain; the cost and time required to
27 conduct systematic reviews; among others. These challenges have been reviewed elsewhere
28 (12, 13).

29 The landscape is further complicated by the increasing interest in setting food and nutrition
30 guidance and policies to lower rates of diet-related chronic diseases, which are a major driver of
31 health care costs in the United States (11). Historically, the Dietary Reference Intakes and the
32 Dietary Guidelines for Americans were established to inform food and nutrient intakes in
33 “apparently healthy” individuals to maintain nutritional adequacy and avoid diseases of nutrient
34 deficiencies. Compared to diet-related chronic diseases, nutritional deficiency in otherwise
35 healthy individuals generally have a single cause, which is a lack of dietary intake of a particular
36 essential nutrient. Furthermore, virtually all healthy individuals respond similarly to dietary
37 deficiency of a particular nutrient in terms of the dose-response relationship and resulting clinical
38 manifestations. This is not the case when diet-related chronic disease is the outcome used for
39 setting food and nutrition guidance, programs, and policies. The etiologies of chronic diseases
40 are highly complex, resulting from the interactions among many essential nutrients and
41 nonessential dietary components. In addition, chronic disease etiologies are modified by
42 differences in individual biology as well as multiple lifestyle factors and exposures including
43 physical activity, sleep, stress, diet, eating behaviours, immune responses, and toxins, among
44 other factors. The contextual factors that modify connections between food and health are even
45 more complex in low- and middle-income country settings. Hence, it is not surprising there is
46 significant population heterogeneity in the diet-chronic disease relationship compared to that

1 between diet and nutrient deficiencies, indicating the need for new approaches to stratify
2 populations to improve the precision of recommendations based on various contexts (14).

3 Population-based diet, food and nutrition recommendations have focused on avoiding essential
4 nutrient deficiencies with consideration for “apparently healthy individuals”, because the disease
5 process can alter nutritional requirements (15). However, when considering chronic disease
6 reduction as an endpoint for nutrient intake recommendations, individuals at risk for or who have
7 chronic disease cannot be excluded because diet-related chronic diseases can initiate as early as
8 during embryonic development and manifest over a lifetime. More than 60% of US adults are
9 affected by a chronic disease, and food- and nutrient-based guidance based on avoidance of
10 nutritional deficiency may not apply to them (16). Globally, essential nutrient deficiencies occur in
11 the obese state. Hence, inclusion of chronic disease outcomes for food and nutrition guidance
12 greatly expands the population under consideration and adds additional heterogeneity in
13 response to dietary and nutrient intake.

14 Consideration of chronic disease endpoints also expands the number of food components under
15 consideration from essential nutrients to any food component that, while not essential, confers a
16 health benefit (17), further increasing the complexity of food and nutrition guideline development.
17 As such, inclusion of chronic disease endpoints in food and nutrition guidance requires expansion
18 of the populations under consideration. Considering these issues, the National Academies of
19 Sciences, Engineering and Medicine recently expanded the definition of the target population for
20 Dietary Reference Intake values to include those with or at risk for chronic disease, with each
21 expert committee being responsible for establishing exceptions that apply specifically to the
22 nutrient(s) under review (18). This expansion of the population under consideration adds to the
23 complexity of data required for establishing recommendations.

24 *Technology Terms and Applications to Nutrition Evidence Synthesis*

25 Broadly defined, AI refers to technologies capable of mimicking human intelligence, including
26 having the capacity to solve complex problems and inclusive of various terms and types of
27 strategies as it pertains to evidence synthesis (19). Over the past decade, AI emerged as an
28 important technology which may provide decision support, early on with specialized deep learning
29 architectures, and more recently with general, pre-trained large language models (LLMs) (20, 21).

30 Modern LLMs have emerged because of: (1) parameter estimation algorithms that make it
31 possible to train models with billions or trillions of parameters; (2) computing infrastructure such
32 as graphics processing units (GPUs) that make it possible to fit models in days or weeks rather
33 than decades but may be cost prohibitive to most researchers; and (3) internet-scale training data
34 corpuses, enabling an arsenal of applications, some of which are deeply embedded in our
35 everyday lives (22, 23). In food and nutrition, AI is now being utilized to guide more precise and
36 accurate food and nutrition guidance to improve health (24). Data science methods offer the
37 opportunity: 1) to automate and thereby accelerate the process of synthesizing data and
38 generating systematic reviews, saving cost and providing decision makers with up-to-date and
39 comprehensive scientific information to make timely decisions; 2) to provide new types of
40 computed data with respect to the complex etiology of the diet-disease relationship, and 3) to
41 identify and classify variation in individual responses to diet. As such, AI offers a timely and cost-
42 effective avenue to develop the strong evidence base necessary to establish effective
43 nutrition/food interventions that prevent and/or manage chronic disease, including data such as
44 electronic medical records (EMRs), as well as take advantage of new types of personalized data

1 from wearables. However, the quality and sparsity of data currently available for such AI-based
2 analyses limit its utility.

3 *Purpose of the Summary*

4 This perspective summarizes findings from a meeting of experts convened at the National
5 Academy of Sciences, Engineering and Medicine (Tables S1, S2). The purpose of the meeting
6 was to examine the current state and future potential of AI in generating new types of computed
7 data as well as automating the generation of systematic reviews to support evidence-based food
8 and nutrition policy, practice and guidance. Participants included expert computational, data and
9 nutrition scientists, as well as scientists from federal research and regulatory agencies. The
10 conference agenda was organized into two main areas (i.e. “parts”) as outlined below.

11 **Part I: Emerging Sources of Scientific Evidence**

12 Establishing scientific recommendations for chronic disease risk reduction through food and
13 nutrition presents an enormous data challenge. This is due to the complexity of food and food
14 components that individuals are exposed to, the variation in individual response to food and
15 nutrient exposures, the number of chronic diseases that are affected by food, the latency and
16 cumulative effects of nutrition on the progression of diet-related diseases that manifest over a
17 lifetime, among many other factors that have been described elsewhere (17). This complexity and
18 the associated costs limit the generation of high-quality scientific evidence through randomized
19 controlled trials that are most often short in duration due to the funding structure of research. The
20 availability of large EMR databases and related real-world health and exposure data, coupled
21 with advances in AI models that mine and automate the synthesis of these resources, provides
22 additional inputs into causal inference models that may provide a less expensive approach to
23 understand the diet-disease relationship and its inherent individual variation. However, EMRs
24 currently have limited data on dietary intakes, nutritional biomarkers and other relevant variables
25 at present. With all AI models, a key consideration is the nature of the training or input data.
26 Cross sectional data, for example, is limited for making causal claims whereas longitudinal data is
27 logistically challenging to collect, suffers from confounding, but represents the longer latency of
28 nutritional exposures and chronic disease risk. Ultimately, the quality of any synthesis of any data
29 relies on the scientific rigor and data available for (i.e. garbage in, garbage out”). Training AI
30 models to advance evidence synthesis can be coordinated and managed by efforts to collect the
31 optimal combinations of data needed to leverage the potential of and amount of data needed to
32 seed AI models.

33 *Lessons learned from cancer drug response prediction models.* Deep neural networks are
34 actively being deployed in sophisticated models for predicting therapeutic responses in cancer
35 (25). However, two major challenges continue to prevent their integration into broader clinical
36 practice (26). The first is lack of model interpretability. The ability to scrutinize the inner workings
37 of a model is critical to building trustworthy AI tools, especially in high-stakes applications such as
38 precision medicine. Visible neural networks (VNN) enable direct model interpretation by mapping
39 the neural network architecture to hierarchical knowledge graphs of biological components and
40 functions (Figure 1) (27, 28). A recently published VNN predicted palbociclib efficacy in breast
41 cancer treatment; it captured 8 molecular assemblies integrating rare and common mutations in
42 90 genes (29). Another recent publication highlighted 41 assemblies involved in modulating
43 response to common chemotherapies (30). These works serve as illustrative proofs-of-concept to
44 develop robust composite biomarkers.

1 A second challenge is related to generalizability. Drug response prediction models are often
2 trained on preclinical datasets. Transferring information from large preclinical datasets to
3 accurately predict treatment response to smaller patient datasets is particularly challenging, and
4 may require careful causal modelling. For predictive tasks, massive pre-trained networks can
5 adapt to new tasks when provided only a handful of examples; this is called “few-shot learning”
6 and was used to perform Translation of Cellular Response Prediction (31). This approach
7 realized better predictive performance across multiple data types, including tumor cell lines,
8 patient-derived tumour cell cultures, and patient-derived tumor xenografts. An independent
9 reusability study was able to apply this approach to two patient cohorts and demonstrated its
10 superior performance (32).

11 VNN models are not yet common in elucidating the role of dietary exposures and their availability
12 on cellular networks and biomarkers of disease etiology. If well executed and validated, these
13 tools can inform biomarker discovery from basic research for clinical utility. This requires the
14 transfer of information across a series of contexts (e.g. from cell lines to patients, from one patient
15 cohort to another, from large populations to small ones or even individuals) with limited data.
16 Similarly, few-shot learning may be applied to transfer biomarkers across contexts.

17 *Knowledge graphs to reveal the etiology of chronic diseases.*

18 Mining multidimensional patient data that includes endome-type data (e.g., clinical exams,
19 laboratory data, imaging, genetics, etc.) and ectome-type data (e.g., age, demographics,
20 exposures, food, social determinants of health) may allow comprehensive consideration of the
21 risk factors underpinning the etiologies of chronic disease. Extracting trustworthy information from
22 large data sets that is both statistically and biologically meaningful, and that can infer causal
23 factors and their relationships, is yet unrealized (33) but is essential for developing effective
24 interventions that are tailored to the context of an individual’s circumstances. Knowledge graphs
25 are a tool to convert large volumes of new data to information and ultimately actionable
26 knowledge but must come from well-established information and include layers of hierarchical
27 organization, their interactions, and relationships across the continuum, including consideration of
28 biological and social complexity. Such bottom-up approaches interconnect layered networks
29 within and across known biological, social, and other domains. The domains can include genes to
30 proteins to pathways (metabolic, signalling, etc.), to cells, organs, the microbiome and the
31 individual within the social context, including the complexity of spaces and locations related to
32 disease incidence, exposures, and temporal changes that individuals experience.

33 One example knowledge graph is the Scalable Precision Medicine Open Knowledge Engine
34 (SPOKE) (34). It has over 40 million concepts that are connected by over 120 established
35 biologically meaningful relationships gathered from existing knowledge in the scientific literature.
36 SPOKE was built by integrating information from more than 50 public databases and contains
37 experimentally determined information on various biological pathways and their architecture, with
38 every node within a network receiving a weighted score relative to its overall importance
39 explaining to risk or function. SPOKE has recently incorporated more than 1000 food items and
40 their relationships to biochemical compounds as determined by mass spectrometry (34) (Figure
41 2). When SPOKE analysed data from six million EMRs, it led to the identification of the nodes of
42 most importance to Parkinson’s disease. SPOKE retrospectively predicted individuals who would
43 develop the disease three years prior to a diagnosis with 83% accuracy and performed similarly
44 to that of clinical expert predictions (35). While not currently clinical grade/caliber, further
45 development and refinement of SPOKE is expected to support its deployment in medical practice.
46 SPOKE includes more than 10,000 disease states and can be used towards discovery and

1 applications related to food, health, and disease. In nutrition research, SPOKE can be used to
2 generate hypotheses by predicting the immediate biological and long-term health effects of
3 consuming individual nutrients and other bioactive food components through a dietary
4 supplement, the optimal combinations of nutrient intakes, and/or effects of consuming specific
5 foods or dietary patterns.

6 *Predictive modelling and individual responses.* The concept of precision nutrition is founded on
7 the premise that identifiable subgroups of individuals respond differently to nutrients, foods, and
8 dietary patterns when chronic disease endpoints are considered (14, 36). The need for precision
9 nutrition is supported by our understanding of human evolution. Human responses to food and
10 nutrition have been under a strong selective pressure in the face of increasing genetic diversity
11 through adaptation to local food environments, which differed considerably across the globe.
12 Adaptation to local food environments enabled population expansion, as classically seen with
13 genetic variation that enabled lactase persistence (37). However, the degree of meaningful
14 biological variation among individuals that necessitates more precision in food and nutrition
15 interventions and recommendations for chronic disease risk reduction remains unresolved. To
16 fully establish the need for greater precision in food and nutrition guidance, two critical questions
17 need to be addressed. First, are the differences among individuals clinically meaningful? Second,
18 do we have predictive biomarkers for the diet-health response?

19 Prediction models are widely used to mine massive data sets to explore the complexity
20 underlying the interactions among endogenous biological factors and environmental exposures
21 that define or relate to human health. Importantly, they offer the possibility of identifying causal
22 dietary and other factors and predicting their intervention response (38). Establishing
23 reliable predictive models of intervention responses has proved challenging due to limitations
24 including bias resulting from several sources, such as algorithmic bias (39), data collected for one
25 purpose being used for other purposes, lack of participant diversity, lack of domain expertise in
26 data selection, among others (40). This, and a dearth of success stories, has led to scepticism for
27 identifying biomarkers that are predictive of medical and nutritional intervention responses. For
28 example, AI prediction models of antipsychotic medications trained on RCT data failed to predict
29 patient outcomes when applied to out-of-sample patients, indicating treatment outcomes are not
30 generalizable for schizophrenia, emphasizing the strong modifying effects of an individual's
31 contexts (41).

32 Traditional clinical trials focusing on nutrition and pharmaceuticals are typically designed to
33 determine the average effect of an intervention, which becomes the evidence-base for
34 establishing generalized population-based applications. These trial designs give less attention to
35 the variation in response among individuals and in fact may mask positive or negative outcomes
36 among subgroups of participants. In contrast, N-of-1 trials seek to identify and characterize
37 variation in responses to multiple interventions provided to the same individual, often separated
38 by wash-out periods, and thereby optimize interventions for that individual (42). N-of-1 trials
39 thereby determine which interventions are better suited for individuals with certain characteristics
40 (43). Such studies that seek to identify and quantify variation around an average response, or a
41 more discreet effect revealing overt responders and non-responders, can be cost-efficient, since
42 the statistical power is optimized when the number of observations is maximized on fewer
43 individuals, compared to fewer observations on more individuals. N-of-1 trials have been used in
44 the fields of psychology and education research, but to a lesser extent in nutrition research.

45 Predictive model reliability can be improved by combining sparse real-world data in large samples
46 with more rigorously collected and outcome-focused data, including data collected during clinical

1 trials. This approach can be more efficient than using sparse data on large number of individuals
2 or very costly yet plentiful experimental data on fewer individuals. For example, models built on
3 massive, randomly sampled, sparse, real-world data, such as the UK Biobank and the NIH-
4 funded All Of Us Study, can be strengthened by calibrating with more sophisticated dense, yet
5 costly, empirical data (44), such as derived from aggregated N-of-1 studies. In this light,
6 aggregated N-of-1 trials might be efficient and appropriate vehicles for vetting or testing the
7 predictions of population-based AI/LLM analyses. Thus, if a new AI/LLM based model is designed
8 to determine which individuals are likely to benefit from a nutritional intervention, then more
9 detailed studies of well-chosen data subsets from individuals for whom predictions were made
10 should shed light on their veracity and expose limitations. This approach is essential to advance
11 the concept of precision nutrition.

12 Other strategies have been employed to strengthen real-world data to understand variation in
13 response (45). AI techniques used to identify factors that are associated with an intervention
14 response are limited by the data sets that they are trained on and cannot be used to infer
15 causation. Training models on more detailed experimental trial data, with limited training on
16 readily available contextual real-world data (e.g. EMRs; large epidemiological datasets),
17 enhances their ability to identify predictive factors and account for variation in individual
18 responses. Such approaches, carefully deployed, have the potential to be more cost-effective and
19 potentially more reliable than conducting large randomized controlled trials.

20 Use of digital twins can also improve predictive models by accounting for the factors that lead to
21 variation in responses. This is achieved by limiting training sets to specified subsets of individuals
22 within a data set who share similar characteristics. Digital twins may better anticipate the
23 health trajectory of a target individual (i.e., 'digital twins' of the target individual) as opposed to
24 using all individuals in the large data set when making predictions about the target individual's
25 health trajectory. Digital twins may share similar genetic, demographic, microbiome, and other
26 characteristics (46-48).

27 *Addressing the complexity of food systems, diets, and their relationship to health.* Food systems,
28 diets, nutrition, and human health exist along a continuum. Dietary patterns differ by context
29 across geography, culture, and socioeconomic status, among other factors (49). Food
30 consumption also has temporal, hedonic, religious, and social dimensions, all of which may relate
31 to health outcomes (14). This has motivated interest in applying AI tools to establish connections
32 across the food value chain and thereby identify opportunities to improve the health-promoting
33 properties of the food system. Traditionally, meta-analyses have been instrumental in
34 understanding the impact of dietary practices and help inform medical decisions. Data science
35 technologies permit a comprehensive approach to addressing food systems and health within
36 these contexts.

37 One example of a dietary pattern that is used clinically is FODMAP. A recent AI model used data
38 from various studies to correlate the success of low FODMAP diets (Fermentable
39 Oligosaccharides, Disaccharides, Monosaccharides, And Polyols) for the treatment of patients
40 with irritable bowel syndrome (IBS), where only 50-70% of the patients respond well to this
41 standard of care treatment. The AI model combined metagenomics and machine learning
42 analysis and provided hypotheses about the mechanism explaining patient segmentation,
43 predicted patient response, and informed treatment decision based on 3 biomarkers (50).
44 Expanding this approach for management of other chronic diseases through diet is an active area
45 of investigation.

1 Knowledge graphs are also playing a key role in assembling and structuring data related to food
2 composition. Agricultural food products contain tens of thousands of chemicals. The FoodData
3 Central database from USDA curates compositional information from 236 foods and 400
4 chemicals that have been validated (51). Recently, there have been advances in streamlining the
5 generation of Knowledge Graphs with using deep Natural Language Processing (NLP)
6 techniques and LLMs to support decision support and accelerated discovery (52). Food Atlas (53)
7 is an AI-generated knowledge graph that has extracted more than 230K food-chemical
8 composition relationships from more than 155K scientific papers, and ranked the confidence level
9 of each relationship based on the existing published evidence (54). This analysis estimated that
10 approximately half of the identified relationships were not previously discovered. While false
11 discovery rate is always a caveat to consider, this lends credence to the potential for utilizing
12 such techniques for discovery. By applying knowledge graph completion methods, new
13 hypotheses can be formed and experimentally validated, providing a framework for automated
14 hypothesis generation. The next version of Food Atlas that is under release, uses a combination
15 of LLMs and hybrid Knowledge Graph Language Models to integrate food, ingredients,
16 chemicals, flavors and health effects.

17 **Part II: Accelerating the Process of Evidence Synthesis**

18 The body of unstructured biomedical data is vast and growing rapidly, hindering physicians' and
19 policymakers' ability to make the most informed decisions grounded in the totality of the evidence
20 base. SR and evidence synthesis are key to developing evidence informed decisions whether
21 they are from a medical, research, or policy lens. However, the process of conducting and
22 publishing systematic reviews is time consuming and expensive, and many of the tasks are highly
23 repetitive but cannot be automated trivially. Consequently, only half of high-quality reviews in
24 biomedical and allied health fields are completed within two years of protocol publication (55).
25 SRs can be expensive to produce and can quickly become outdated, sometimes even before
26 they are published (56), lending credence for the need for newer methods that function in real-
27 time. Study screening, data extraction and synthesis are key bottlenecks in generating systematic
28 reviews. There is a need to design, implement, and deploy NLP tasks, corpora, and models to
29 help domain experts navigate and make sense of the vast array of biomedical evidence, ranging
30 from notes in EMRs to published reports of clinical trials, which are generally stored as
31 unstructured text and therefore not readily accessible or mineable.

32 High quality evidence synthesis adheres to the principles of transparency, reproducibility, and
33 methodological rigor, following prespecified processes (57, 58). Otherwise, SR
34 findings/conclusions can be highly dependent or influenced by subjective judgements (59, 60). It
35 is these and related challenges that motivated the development of AI tools for SR, but the uptake
36 has been slow (61). By necessity and logic, the process must include human judgement or
37 oversight in the identification of the relevant literature base from raw search results (based on
38 pre-specified search criteria that is then screened) as well as in the rating the risk of bias of
39 individual studies and grading the overall certainty of the available evidence (3). Literature
40 screening, usually conducted manually by human non-content experts (e.g. trainees, students,
41 contractors), is the most time and resource intensive stage of the process and can be subject to
42 various types of bias. To mitigate bias and the temporal currency of the SR process, human-AI
43 hybrid approaches have been developed, and evaluated for their effectiveness, in accelerating
44 the generation of high-quality evidence synthesis products that promote timely evidence-based
45 scientific guidance for decision makers.

1 Goals of including AI applications in the evidence synthesis process include accelerating
2 innovation and time-to-completion, improving productivity, and cost reduction (62). Title and
3 abstract screening for inclusion in a SR generally reduces the number of studies identified
4 through a literature search by 95%, and hence is a task that is well suited for automated text
5 classification. Early NLP models were frequency-based models, classifying studies by the
6 frequency of individual terms within a document/text. More recent approaches use neural
7 network-based methods, up to and including LLMs. Currently, there are several AI-powered
8 screening tools (both commercial and open-source) available to accelerate title and abstract
9 screening but rely on frequency-based models (e.g., EPPI-Reviewer, abstrackr, DistillerSR,
10 RobotReviewer and Rayyan) that represent the industry standard (63-66). As an example, the
11 USDA Nutrition Evidence Systematic Review group, which conducts systematic reviews in
12 support of establishing the Dietary Guidelines for Americans, uses AI-powered screening tools
13 (3).

14 An early and relatively large pre-trained neural network was the Bidirectional Encoder
15 Representations from Transformers (BERT). BERT is pre-trained on a large volume of text, and
16 can be fine-tuned for particular tasks. This model has been incorporated into human-AI hybrid
17 evidence synthesis teams (62). The collaborative screening process involves subject matter
18 experts identifying the screening criteria, followed by the training of the NLP using a limited
19 number of studies screened by subject matter experts. Once the model is judged to function
20 adequately, it ranks new documents never seen by the model (62). A final review of all selected
21 documents is conducted by experts. The approach is iterative as feedback from the experts
22 continuously constrains and improves the model. The approach may incorporate active learning
23 into the human-AI hybrid team by exploring and testing different sampling strategies, including
24 random sampling, least confidence sampling, and highest priority sampling, and evaluating their
25 effectiveness on the collaborative screening process.

26 Incorporating the BERT-based AI agent into a human team was found to reduce the human
27 screening effort, including the number of documents that humans need to read, by 68.5%
28 compared to the case of no AI assistance, and by 16.8% compared to the industry standard that
29 uses a frequency-based language model and a support vector machine-based classifier (Figure
30 3). These values are for the human screening effort required to identify 80% of all relevant
31 documents. The process was further improved by applying a HP sampling strategy to the human
32 screening effort, resulting in 78.3% reduction in human screening effort to identify 80% of all
33 relevant documents compared to no AI assistance. The BERT-based model uniformly
34 outperformed the industry standard NLPs in classification performance.

35 Key limitations to using active learning-enhanced human-AI hybrid team workflow process are the
36 time of communication among subject matter experts and computational scientists, the level of
37 measurement error inherent to human labels, which is addressed with additional iterative training,
38 and trust among the experts and the model. Future expansion to full text screening is expected to
39 improve classification performance but can be limited by inaccessible documents that are not
40 published in open-access format. It is important to note that the field of LLMs is changing rapidly
41 and becoming more powerful with generative models which could improve accuracy and be able
42 to summarize evidence but will require validation.

43 *Extracting and synthesizing medical evidence with LLMs.* Clinical trial results are disseminated
44 through natural language articles and hence are largely unstructured or semi-structured, including
45 clinical trial databases such as clinicaltrials.gov. NLP methods in general, and automated
46 summarization in particular, offer a potential means of helping domain experts identify and make

1 better use of the totality of scientific data to inform treatment and other-related decisions. Variants
2 of LLMs are being used to extract and structure findings from clinical trial reports, and to generate
3 automatic summaries of all published evidence pertaining to a particular clinical question. An
4 available prototype, Trialstreamer, is a publicly available living repository of all articles describing
5 RCTs in humans that makes RCT data fully computable (64, 67) (Figure 4). It monitors PubMed
6 and other sources daily, then structures the data using models that extract and tabulate key
7 information including PICO (Population, Intervention, Comparison, Outcome) element information
8 and other metrics such as sample sizes. Trialstreamer can conduct aspects of Cochrane-style
9 risk of bias assessments, such as whether a trial was randomized or blinded, which otherwise
10 involves subjective judgements by humans. Trialstreamer can infer main findings of a study
11 through a semi-automated process that accelerates human assessment by about 30% (68), and
12 the results are generally in agreement with human assessments. The database can be searched
13 for all studies relevant to a well-formed clinical question if indexed by PubMed (emerging pre-
14 publication websites, by lack of per-review, are not incorporated).

15 In development for the next iteration of Trialstreamer is the capability to generate Cochrane-style
16 systematic reviews, including meta-analyses, and a natural language narrative that describes the
17 summary of results. Current technologies permit automatic generation of plausible summaries but
18 may, or even often, include “hallucinations” in the conclusions which is a real problem that needs
19 to be addressed to ensure “trustworthy” information. Other limitations pertain to the assessment
20 of more nuanced information from studies, such as extraction and critical appraisal of intervention
21 and outcome ascertainment methods given the discipline- and method-specific nature of this kind
22 of data.

23 *Key performance indicators for AI-assisted evidence synthesis.* Looking forward, automated
24 evidence synthesis products must be fit for purpose, and the evidence synthesis processes
25 should be robust a predictably changing environment (e.g., the increasing rate at which primary
26 research is published) and rapidly responsive to unpredictable shocks (e.g., health emergencies
27 such as the COVID-19 pandemic). This will require new tools and processes, but should also
28 build upon an understanding of three key performance indicators (KPIs): 1) time use and time to
29 completion; 2) resource use and economic sustainability, and 3) correctness (69, 70). Shaping
30 the future of evidence synthesis, both technologically and culturally, is essential to ensure that it
31 continues to meet stakeholder needs.

32 The three KPIs have been assessed in a limited number of cases. A study by Tercero-Hidalgo
33 examined the influence of using AI in the systematic review process related to COVID-19 (71).
34 The prespecified study included 3,999 systematic reviews, 28 of which used AI. The use of AI
35 was associated with publication in journals with a higher impact factor (8.9 vs. 3.5), more
36 abstracts screened per author (302 vs. 140) and fewer texts screened per author (5.3 vs. 14 full
37 texts) but curiously no effect on time to completion. In another prespecified study, Meneses -
38 Echavez et al. examined the KPIs person hours and time-to completion prior to and following
39 adoption of ML in the systematic review process from August 2020 to January 2023 at the
40 *Norwegian Institute of Public Health* (70). This study also found using ML required more person-
41 hours and other resources, with no effect on time to completion.

42 The third KPI, correctness, is the most difficult to assess, but could be evaluated by: 1) comparing
43 AI outputs to human reviewers, who are assumed to be making correct decisions; 2) comparing
44 AI outputs to results such as meta-analytical estimates from closed reviews under the assumption
45 that findings in closed reviews are sufficiently close to the truth (reviews are closed if adding
46 additional studies is expected not to change the existing findings); and 3) a simulation approach

1 in which AI tools for evidence synthesis are applied to bodies of literature using computed data,
2 generated using models such as LLMs, where the true values of effect measures such as hazard
3 ratios are known by construction, facilitating comparison of AI outputs with known ground truth.
4 To date, only the first approach has been used for analyses, and is biased by the assumption that
5 human reviewers are correct when in reality they can introduce inconsistency due to human
6 judgement. Vembye and Dietrichson (unpublished data) compared the performance of non-expert
7 reviewers (students) compared to the GPT-4 model for title and abstract screening using results
8 of a literature search that yielded 4,136 articles. The GTP-4 model achieved 90% sensitivity and
9 94% specificity, and where nonexpert humans and the LLM disagreed, subject matter expert
10 humans generally agreed with the LLM.

11 Looking forward, there are many limitations of AI approaches that must be overcome to achieve
12 correctness. It is recognized that there is a trade-off between accuracy and confidence with time
13 savings and efficiency when automating evidence synthesis. Understanding what type of
14 scientific product is needed for a particular purpose (e.g. guideline development) where the need
15 for comprehensiveness and accuracy versus expediency can be pre-specified and reported
16 transparently, otherwise cheap, fast, and possibly incorrect evidence synthesis may result.

17 AI tools may also be abused to quickly produce poor-quality “reviews”, which poses new threats
18 to evidence synthesis. LLMs may also facilitate the production of fake, fraudulent, or flawed
19 primary studies (e.g., zombie trials). It is estimated that hundreds of thousands of zombie trials
20 already circulate in the literature, and their inclusion in evidence syntheses is problematic (72).
21 Furthermore, online AI tools are vulnerable to digital attack including denial-of-service attacks and
22 data set poisoning (73). Other concerns include privacy violations, underrepresentation of studies
23 in minority languages, and the commercial interests of companies marketing AI tools out of
24 alignment with stakeholder needs.

25 Finally, AI tools are perhaps only necessary because scientific results are not reported using
26 standardized structured data formats that permit accurate and comprehensive automated search
27 and data extraction across the entire literature. While reports for some trials are available in
28 machine-readable formats such as JSON and FHIR from clinicaltrials.gov, future work could
29 focus on dramatically extending the coverage and depth of scientific reporting, perhaps using
30 fine-grained and federated graph databases and standardized ontologies.

31 **Discussion**

32 Advances in AI are providing decision makers new ways of accessing and making sense of
33 scientific evidence. Although AI tools alone cannot generate evidence *de novo*, they are capable
34 of processing, daisy-chaining and/or merging evidence across existing datasets into new formats.
35 They have been used to create synthetic dose-response relationships drawing on pathway data
36 from different data sets, which have aided authoritative organizations in setting food and nutrition
37 policy (74, 75). However, the trustworthiness of computed data, including information from VNN
38 and knowledge graphs, and its relative positioning in the hierarchy of evidence has not been
39 addressed (76).

40 The established evidence hierarchy describes the strength of data types based on study design
41 as they relate to causal inference. As one moves up the hierarchy it is assumed that study quality
42 increases and risk of bias decreases, and thereby the certainty of relationships between
43 interventions/exposures and outcomes is higher (76). Well-designed Randomized Controlled
44 Trials, which sit at the top of the hierarchy, can determine causal relationships. As such,

1 systematic reviews, and meta-analyses of these trials are considered the highest level of
2 evidence (76). However, like traditional Randomized Controlled Trial designs, systematic reviews
3 and meta-analyses generally emphasize average responses across many studies, and often fail
4 to consider variation in response between studies or individuals (77). VNNs and knowledge
5 graphs provide the opportunity to address overall effects of an intervention, as well as address
6 variation in response among individuals, but their potential to determine causality has not been
7 established (78), nor has there been consideration to how computed evidence compares to other
8 traditional types of evidence. The quality of AI-assisted SR is also dependent on the body of
9 literature available.

10 Limitations to the established hierarchy-of-evidence include uncertain generalizability of the
11 findings, even when the evidence for causation is strong. The lack of generalizability is rooted in
12 biological heterogeneity within populations that contributes to variance in the exposure-outcome
13 relationship. Likewise, social, environmental, and other contexts in free-living populations can
14 influence efficacy and effectiveness of interventions or exposures. These effects on context limit
15 the ability to predict nutrition intervention outcomes in low- and middle-income countries based on
16 relationships and contexts established in high income countries. Furthermore, the strength of
17 evidence does not always inform whether interventions will have a meaningful magnitude of effect
18 that has a clinical and public health value even when causal inference is strong. Knowledge
19 graphs consider the many biological and social dimensions of food, individuals and health. Their
20 application to nutrition questions, especially when combined with LLMs, presents an exciting and
21 transformational opportunity to connect food and health in a way that considers individuals and
22 their contexts.

23 Ideally, computed data will lead to multiple new types of evidence that will be available to
24 decisionmakers, yet frameworks and appraisal tools do not exist to guide their use. Rather than a
25 single hierarchy, there is an increasing need for a multi-dimensional assessment of the totality of
26 the evidence that is fit for purpose, considers the properties of the evidence and how the
27 outcomes are affected in multidimensional situations. Such a framework should consider and
28 potentially rank the properties of different forms of scientific evidence including causality,
29 generalizability, risk of bias, precision, dose-response, and magnitude of effect, and their relative
30 importance for different purposes.

31 Decisionmakers emphasize the need to accelerate the synthesis of scientific data in response to
32 emergent and sustained societal needs. This includes outcomes of efficacy, effectiveness, and
33 equity across a population. Understanding the generalizability of even the strongest scientific
34 evidence is also essential, as many policy decisions are made locally and include contextual
35 realities in which research and policy making is done. Automating the SR process and
36 incorporating computed evidence can address many of these concerns. For example, elements of
37 equity can be improved by including data reported in underrepresented languages, which are
38 often excluded, through LLMs.

39 In the ideal case, automated real-time collection and analyses of data of high relevance to clinical
40 and public health from all sources is the goal. This will allow more rapid science-informed policies
41 and create a continuously learning health system. Learning systems characterized by automated
42 real-time collection and analyses of data in nutrition could facilitate regular updates to both the
43 DGAs and DRIs as new data become available through a semi-automated process that includes
44 expert input and review (79).

1 AI can also inform future research priorities. AI approaches can assist research funding agencies
2 in identifying gaps in knowledge (identify holes or uncertainty in networks) in real time to guide
3 and prioritize high impact research needs that have a high societal return on investment,
4 especially concerning both continuing and emerging public health threats, including setting
5 priorities for the Dietary Guidelines for Americans.

6 Trust in food and nutrition research is essential, otherwise science-informed guidance and
7 recommendations will not achieve or will diminish the impact of their intended health outcomes
8 (80). The inclusion of validated and reliable data science tools into the process of food and
9 nutrition research, and its translation for public benefit, offers the opportunity to increase public
10 trust. This will be challenging as current LLMs and other tools are essentially “black boxes”; no
11 one knows exactly how they work, or when they will “hallucinate” rather than provide correct
12 information. While this may be less of a concern when these technologies are used in an
13 analytical mode to screen, identify, or extract straightforward evidence during semiautomated
14 evidence synthesis, applications of the technology that generate computed evidence will have to
15 be carefully validated, replicated, and communicated transparently. On the other hand, data
16 science tools offer the potential for more personalized nutrition guidance where individuals can
17 access the science and realize the benefit, as opposed to generalized recommendations that may
18 not be optimal for everyone. These tools also offer the opportunity to reduce bias in nutrition.
19 While data from individuals of European ancestry is overrepresented relative to US demographics
20 in many health-related databases, AI tools such as digital twin approaches may allow us to
21 minimize or eliminate biases and data misalignments by moving away from population averages
22 that might poorly reflect underrepresented individuals and towards causal inferences and
23 predictions that address the unique characteristics of individuals.

24 Finally, meaningful advances in the application of AI to nutrition research, policy and practice will
25 require the inclusion of more, consistently collected, richer nutrition and diet data in EMRs,
26 greater engagement of data scientists with nutrition scientists, and ensuring the next generation
27 of nutrition scientists are trained in the data sciences.

28 **Acknowledgments**

29 We are grateful to the Food and Nutrition Board of the National Academies of Sciences,
30 Engineering and Medicine for hosting the meeting of experts. We are grateful to Mikkel Holding
31 Vembye and Jens Dietrichson at The Danish Center for Social Science Research, Aarhus and
32 Copenhagen, Denmark, for sharing their title and abstract screening results. We wish to
33 acknowledge the attendance and contributions of meetings participants listed in **Supplemental**
34 **Table 1**.

35 **Funding:**

36 This work was funded by a grant from the Bill and Melinda Gates Foundation, MN-Systematic
37 approach to evaluate nutrition biomarkers for MNCH outcomes - INV-047386.
38

39 **Data Availability Statement:**

40 There are no data underlying this work.
41

42 **References**

- 43 1. G. A. Kelley, K. S. Kelley, Systematic reviews and meta-analysis in nutrition research. *Br J Nutr* **122**, 1279-1294 (2019).
- 44 2. P. M. Brannon, C. L. Taylor, P. M. Coates, Use and applications of systematic reviews in
45 public health nutrition. *Annu Rev Nutr* **34**, 401-419 (2014).
- 46
- 47
- 48

- 1 3. M. K. Spill *et al.*, Perspective: USDA Nutrition Evidence Systematic Review Methodology: Grading the Strength of Evidence in Nutrition- and Public Health-Related Systematic Reviews. *Adv Nutr* **13**, 982-991 (2022).
- 2
- 3
- 4 4. P. R. Trumbo, S. I. Barr, S. P. Murphy, A. A. Yates, Dietary reference intakes: cases of appropriate and inappropriate uses. *Nutr Rev* **71**, 657-664 (2013).
- 5
- 6 5. S. P. Murphy, Using DRIs as the basis for dietary guidelines. *Asia Pac J Clin Nutr* **17 Suppl** **1**, 52-54 (2008).
- 7
- 8 6. 2003)in *Dietary Reference Intakes: Guiding Principles for Nutrition Labeling and Fortification* (Washington (DC)).
- 9
- 10 7. S. P. Murphy, A. A. Yates, S. A. Atkinson, S. I. Barr, J. Dwyer, History of Nutrition: The Long Road Leading to the Dietary Reference Intakes for the United States and Canada. *Adv Nutr* **7**, 157-168 (2016).
- 11
- 12
- 13 8. 1989)in *Recommended Dietary Allowances: 10th Edition* (Washington (DC)).
- 14 9. 2011)in *Child and Adult Care Food Program: Aligning Dietary Guidance for All*, eds S. P. Murphy, A. L. Yaktine, C. West Sutor, S. Moats (Washington (DC)).
- 15
- 16 10. E. R. Service (2023) What is agriculture's share of the overall U.S. economy?
- 17 11. T. V. Jardim *et al.*, Cardiometabolic disease costs associated with suboptimal diet in the United States: A cost analysis based on a microsimulation model. *PLoS Med* **16**, e1002981 (2019).
- 18
- 19
- 20 12. P. M. Brannon *et al.*, Scanning for new evidence to prioritize updates to the Dietary Reference Intakes: case studies for thiamin and phosphorus. *Am J Clin Nutr* **104**, 1366-1377 (2016).
- 21
- 22
- 23 13. M. S. Field *et al.*, Scanning the evidence: process and lessons learned from an evidence scan of riboflavin to inform decisions on updating the riboflavin dietary reference intakes. *Am J Clin Nutr* **116**, 299-302 (2022).
- 24
- 25
- 26 14. R. L. Bailey, P. J. Stover, Precision Nutrition: The Hype Is Exceeding the Science and Evidentiary Standards Needed to Inform Public Health Recommendations for Prevention of Chronic Disease. *Annu Rev Nutr* **43**, 385-407 (2023).
- 27
- 28
- 29 15. P. J. Stover, C. Garza, J. Durga, M. S. Field, Emerging Concepts in Nutrient Needs. *J Nutr* **150**, 2593S-2601S (2020).
- 30
- 31 16. 2022) Chronic Disease Fact Sheets. (National Center for Chronic Disease Prevention and Health Promotion).
- 32
- 33 17. E. A. Yetley *et al.*, Options for basing Dietary Reference Intakes (DRIs) on chronic disease endpoints: report from a joint US-/Canadian-sponsored working group. *Am J Clin Nutr* **105**, 249S-285S (2017).
- 34
- 35
- 36 18. E. National Academies of Sciences, and Medicine. (2022) Defining Populations for Dietary Reference Intake Recommendations: A Letter Report. (The National Academies Press, Washington, DC).
- 37
- 38
- 39 19. J. M. Helm *et al.*, Machine Learning and Artificial Intelligence: Definitions, Applications, and Future Directions. *Curr Rev Musculoskelet Med* **13**, 69-76 (2020).
- 40
- 41 20. A. J. Thirunavukarasu *et al.*, Large language models in medicine. *Nat Med* **29**, 1930-1940 (2023).
- 42
- 43 21. E. Guo *et al.*, Automated Paper Screening for Clinical Reviews Using Large Language Models: Data Analysis Study. *J Med Internet Res* **26**, e48996 (2024).
- 44
- 45 22. A. Vaswani *et al.*, Attention is All you Need. *Advances in Neural Information Processing Systems* **30** (2017).
- 46

- 1 23. T. Wu *et al.*, A Brief Overview of ChatGPT: The History, Status Quo and Potential Future
2 Development. *IEEE/CAA Journal of Automatica Sinica* **10**, 1122-1136 (2023).
- 3 24. A. Eetemadi *et al.*, The Computational Diet: A Review of Computational Methods Across
4 Diet, Microbiome, and Health. *Front Microbiol* **11**, 393 (2020).
- 5 25. A. Partin *et al.*, Deep learning methods for drug response prediction in cancer:
6 Predominant and emerging trends. *Front Med (Lausanne)* **10**, 1086097 (2023).
- 7 26. R. Roscher, B. Bohn, M. F. Duarte, J. Garcke, Explainable machine learning for scientific
8 insights and discoveries. *IEEE Access* **8** (20220).
- 9 27. J. Ma *et al.*, Using deep learning to model the hierarchical structure and function of a
10 cell. *Nat Methods* **15**, 290-298 (2018).
- 11 28. B. M. Kuenzi *et al.*, Predicting Drug Response and Synergy Using a Deep Learning Model
12 of Human Cancer Cells. *Cancer Cell* **38**, 672-684 e676 (2020).
- 13 29. S. Park *et al.*, A deep learning model of tumor cell architecture elucidates response and
14 resistance to CDK4/6 inhibitors. *Nat Cancer* 10.1038/s43018-024-00740-1 (2024).
- 15 30. X. Zhao *et al.*, Cancer Mutations Converge on a Collection of Protein Assemblies to
16 Predict Resistance to Replication Stress. *Cancer Discov* **14**, 508-523 (2024).
- 17 31. J. Ma *et al.*, Few-shot learning creates predictive models of drug response that translate
18 from high-throughput screens to individual patients. *Nat Cancer* **2**, 233-244 (2021).
- 19 32. E. So, F. Yu, B. Wang, B. Haibe-Kains, Reusability report: Evaluating reproducibility and
20 reusability of a fine-tuned model to predict drug response in cancer patient samples.
21 *Nature Machine Intelligence* **5**, 792-798 (2023).
- 22 33. M. Z. Naser, Causality and causal inference for engineers: Beyond correlation,
23 regression, prediction and artificial intelligence. *WIREs Data Mining and Knowledge*
24 *Discovery* 10.1002/widm.1533 (2024).
- 25 34. J. H. Morris *et al.*, The scalable precision medicine open knowledge engine (SPOKE): a
26 massive knowledge graph of biomedical information. *Bioinformatics* **39** (2023).
- 27 35. K. Soman *et al.*, Early detection of Parkinson's disease through enriching the electronic
28 health record using a biomedical knowledge graph. *Front Med (Lausanne)* **10**, 1081087
29 (2023).
- 30 36. in Guiding Principles for Developing Dietary Reference Intakes Based on Chronic
31 Disease, M. P. Oria, S. Kumanyika, Eds. (Washington (DC), 2017), 10.17226/24828.
- 32 37. W. P. T. James *et al.*, Nutrition and its role in human evolution. *J Intern Med* **285**, 533-
33 549 (2019).
- 34 38. I. J. Dahabreh, K. Bibbins-Domingo, Causal Inference About the Effects of Interventions
35 From Observational Studies in Medical Journals. *JAMA* 10.1001/jama.2024.7741 (2024).
- 36 39. L. Messeri, M. J. Crockett, Artificial intelligence and illusions of understanding in
37 scientific research. *Nature* **627**, 49-58 (2024).
- 38 40. A. S. Lea, D. S. Jones, Mind the Gap - Machine Learning, Dataset Shift, and History in the
39 Age of Clinical Algorithms. *N Engl J Med* **390**, 293-295 (2024).
- 40 41. A. M. Chekroud *et al.*, Illusory generalizability of clinical prediction models. *Science* **383**,
41 164-167 (2024).
- 42 42. N. J. Schork, B. Beaulieu-Jones, W. S. Liang, S. Smalley, L. H. Goetz, Exploring human
43 biology with N-of-1 clinical trials. *Camb Prism Precis Med* **1** (2023).
- 44 43. T. Potter, R. Vieira, B. de Roos, Perspective: Application of N-of-1 Methods in
45 Personalized Nutrition Research. *Adv Nutr* **12**, 579-589 (2021).
- 46 44. A. N. Angelopoulos, S. Bates, C. Fannjiang, M. I. Jordan, T. Zrnica, Prediction-powered
47 inference. *Science* **382**, 669-674 (2023).

- 1 45. P. Shah *et al.*, Artificial intelligence and machine learning in clinical development: a
2 translational perspective. *NPJ Digit Med* **2**, 69 (2019).
- 3 46. E. Katsoulakis *et al.*, Digital twins for health: a scoping review. *npj Digital Medicine* **7**
4 (2024).
- 5 47. T. Sun, X. He, Z. Li, Digital twin in healthcare: Recent updates and challenges. *Digit*
6 *Health* **9**, 20552076221149651 (2023).
- 7 48. 2023) Opportunities and Challenges for Digital Twins in Biomedical Research:
8 Proceedings of a Workshop—in Brief. ed L. Casola (National Academies of Sciences,
9 Engineering, and Medicine; National Academy of Engineering; Division on Earth and Life
10 Studies; Division on Engineering and Physical Sciences; Board on Life Sciences; Board on
11 Atmospheric Sciences and Climate; Computer Science and Telecommunications Board;
12 Board on Mathematical Sciences and Analytics, Washington (DC)).
- 13 49. J. Reedy *et al.*, Evaluation of the Healthy Eating Index-2015. *J Acad Nutr Diet* **118**, 1622-
14 1633 (2018).
- 15 50. A. Eetemadi, I. Tagkopoulos, Methane and fatty acid metabolism pathways are
16 predictive of Low-FODMAP diet efficacy for patients with irritable bowel syndrome. *Clin*
17 *Nutr* **40**, 4414-4421 (2021).
- 18 51. E. M. Jennings-Dobbs, S. M. Forester, A. Drewnowski, Visualizing Data Interoperability
19 for Food Systems Sustainability Research-From Spider Webs to Neural Networks. *Curr*
20 *Dev Nutr* **7**, 102006 (2023).
- 21 52. J. Youn, N. Rai, I. Tagkopoulos, Knowledge integration and decision support for
22 accelerated discovery of antibiotic resistance genes. *Nat Commun* **13**, 2360 (2022).
- 23 53. 2024) FoodAtlas. (The AI Institute for Next Generation Food Systems, University of
24 California at Davis).
- 25 54. J. Youn, F. Li, G. Simmons, S. Kim, I. Tagkopoulos, FoodAtlas: Automated Knowledge
26 Extraction 1 of Food and Chemicals from Literature. *bioRxiv* (2024).
- 27 55. M. Z. Andersen, S. Gulen, S. Fønnes, K. Andresen, J. Rosenberg, Half of Cochrane reviews
28 were published more than 2 years after the protocol. *J Clin Epidemiol* **124**, 85-93 (2020).
- 29 56. K. G. Shojania *et al.*, How quickly do systematic reviews go out of date? A survival
30 analysis. *Ann Intern Med* **147**, 224-233 (2007).
- 31 57. A. D. Oxman, G. H. Guyatt, The science of reviewing research. *Ann N Y Acad Sci* **703**, 125-
32 133; discussion 133-124 (1993).
- 33 58. M. Cumpston *et al.*, Updated guidance for trusted systematic reviews: a new edition of
34 the Cochrane Handbook for Systematic Reviews of Interventions. *Cochrane Database*
35 *Syst Rev* **10**, ED000142 (2019).
- 36 59. N. Konsgen *et al.*, Inter-review agreement of risk-of-bias judgments varied in Cochrane
37 reviews. *J Clin Epidemiol* **120**, 25-32 (2020).
- 38 60. S. Minozzi, M. Cinquini, S. Gianola, M. Gonzalez-Lorenzo, R. Banzi, The revised Cochrane
39 risk of bias tool for randomized trials (RoB 2) showed low interrater reliability and
40 challenges in its application. *J Clin Epidemiol* **126**, 37-44 (2020).
- 41 61. A. M. O'Connor *et al.*, Still moving toward automation of the systematic review process:
42 a summary of discussions at the third meeting of the International Collaboration for
43 Automation of Systematic Reviews (ICASR). *Syst Rev* **8**, 57 (2019).
- 44 62. K. M. Edwards *et al.*, ADVISE: Accelerating the Creation of Evidence Synthesis for Global
45 Development Using Natural Language Processing-Supported Human Artificial
46 Intelligence Collaboration. *J. Mech. Des.* **146**, 14 (2024).

- 1 63. I. Shemilt *et al.*, Cost-effectiveness of Microsoft Academic Graph with machine learning
2 for automated study identification in a living map of coronavirus disease 2019 (COVID-
3 19) research. *Wellcome Open Res* **6**, 210 (2021).
- 4 64. S. Ramprasad, I. J. Marshall, D. J. McInerney, B. C. Wallace, Automatically Summarizing
5 Evidence from Clinical Trials: A Prototype Highlighting Current Challenges. *Proc Conf*
6 *Assoc Comput Linguist Meet* **2023**, 236-247 (2023).
- 7 65. M. Ouzzani, H. Hammady, Z. Fedorowicz, A. Elmagarmid, Rayyan-a web and mobile app
8 for systematic reviews. *Syst Rev* **5**, 210 (2016).
- 9 66. C. J. L. Murray *et al.*, Five insights from the Global Burden of Disease Study 2019. *The*
10 *Lancet* **396**, 1135-1159 (2020).
- 11 67. I. J. Marshall *et al.*, Trialstreamer: A living, automatically updated database of clinical
12 trial reports. *J Am Med Inform Assoc* **27**, 1903-1912 (2020).
- 13 68. F. Soboczenski *et al.*, Machine learning to help researchers evaluate biases in clinical
14 trials: a prospective, randomized user study. *BMC Med Inform Decis Mak* **19**, 96 (2019).
- 15 69. J. Clark, C. McFarlane, G. Cleo, C. Ishikawa Ramos, S. Marshall, The Impact of Systematic
16 Review Automation Tools on Methodological Quality and Time Taken to Complete
17 Systematic Review Tasks: Case Study. *JMIR Med Educ* **7**, e24418 (2021).
- 18 70. A. E. Muller *et al.*, The effect of machine learning tools for evidence synthesis on
19 resource use and time-to-completion: protocol for a retrospective pilot study. *Syst Rev*
20 **12**, 7 (2023).
- 21 71. J. R. Tercero-Hidalgo *et al.*, Artificial intelligence in COVID-19 evidence syntheses was
22 underutilized, but impactful: a methodological study. *J Clin Epidemiol* **148**, 124-134
23 (2022).
- 24 72. J. P. A. Ioannidis, Hundreds of thousands of zombie randomised trials circulate among
25 us. *Anaesthesia* **76**, 444-447 (2021).
- 26 73. N. Carlini *et al.*, Poisoning Web-Scale Training Datasets is Practical. *arXiv:2302.10149*
27 (2024).
- 28 74. K. S. Crider, Y. P. Qi, O. Devine, S. C. Tinker, R. J. Berry, Modeling the impact of folic acid
29 fortification and supplementation on red blood cell folate concentrations and predicted
30 neural tube defect risk in the United States: have we reached optimal prevention? *Am J*
31 *Clin Nutr* **107**, 1027-1034 (2018).
- 32 75. K. S. Crider *et al.*, Population red blood cell folate concentrations for prevention of
33 neural tube defects: Bayesian model. *BMJ* **349**, g4554 (2014).
- 34 76. B. Djulbegovic, G. H. Guyatt, Progress in evidence-based medicine: a quarter century on.
35 *Lancet* **390**, 415-423 (2017).
- 36 77. J. IntHout, J. P. Ioannidis, M. M. Rovers, J. J. Goeman, Plea for routinely presenting
37 prediction intervals in meta-analysis. *BMJ Open* **6**, e010247 (2016).
- 38 78. M. A. Hernan, The C-Word: Scientific Euphemisms Do Not Improve Causal Inference
39 From Observational Data. *Am J Public Health* **108**, 616-619 (2018).
- 40 79. 2024) About Learning Health Systems. (Agency for Healthcare Research and Quality).
- 41 80. C. Garza *et al.*, Best practices in nutrition science to earn and keep the public's trust. *Am*
42 *J Clin Nutr* **109**, 225-243 (2019).
- 43
44

1 **Figure Legends:**

2
3
4 **Figure 1. Visible neural network (VNN).** Adapted from Park, Silva, Singhal, et al. (29).
5 The first layer of the VNN incorporates gene-level features, including gene mutations, copy
6 number amplifications (CNA), and copy number deletions (CND). Subsequent assembly layers
7 aggregate gene-level features into assembly-level information, guided by the hierarchical
8 relationships defined by a map of protein assemblies. The output state of each gene (g) and
9 assembly (O) is represented by artificial neurons (one neuron per gene, multiple neurons per
10 assembly). Each node in the hierarchy indicates a protein assembly. An example path of
11 information flow is shown in red.

12
13 **Figure 2.** Scalable Precision Medicine Open Knowledge Engine (SPOKE). The SPOKE
14 biomedical knowledge graph draws upon and integrates over 45 databases.

15
16 **Figure 3.** A human-AI workflow for document screening in evidence synthesis. In Stage 1,
17 experts specify screening criteria for documents, then screen a subset of the documents-of-
18 interest for inclusion or exclusion in an evidence synthesis product. In Stage 2, an AI model is
19 trained on the expert labels of screened documents, and then performs screening of additional
20 documents. In Stage 3, expert labellers evaluate the AI's screening decisions. The final validated
21 screening decisions are used to iteratively re-train the AI model.

22
23 **Figure 4.** Trialstreamer: A living automated automatically updated database of clinical trial reports
24 (67).

25

ACCEPTED MANUSCRIPT

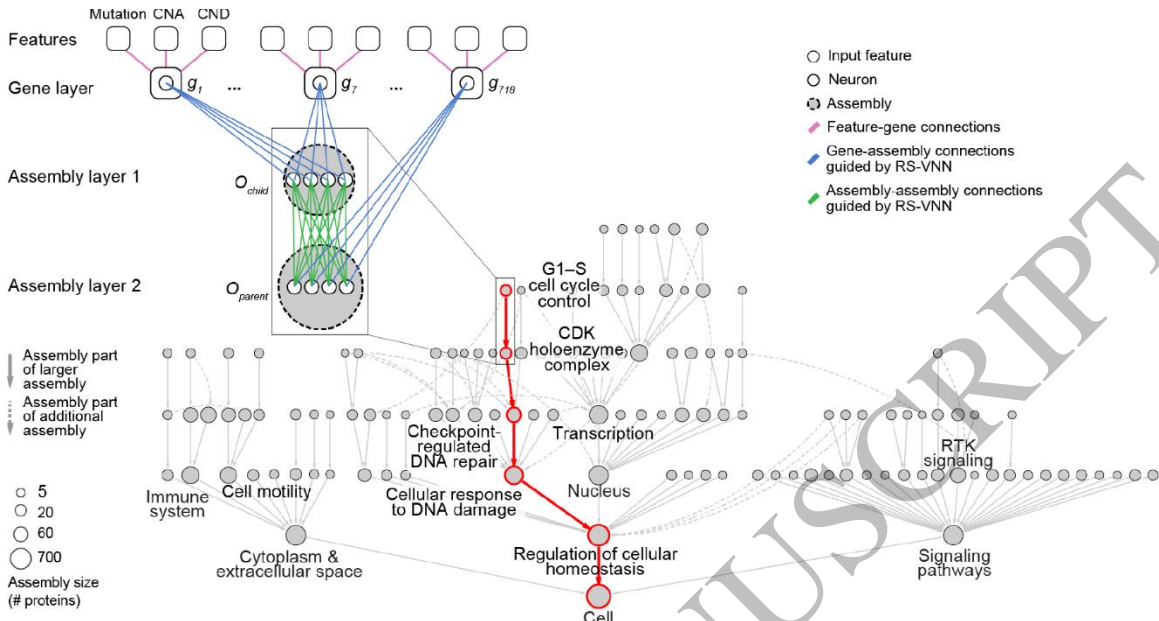


Figure 1
152x81 mm (x DPI)

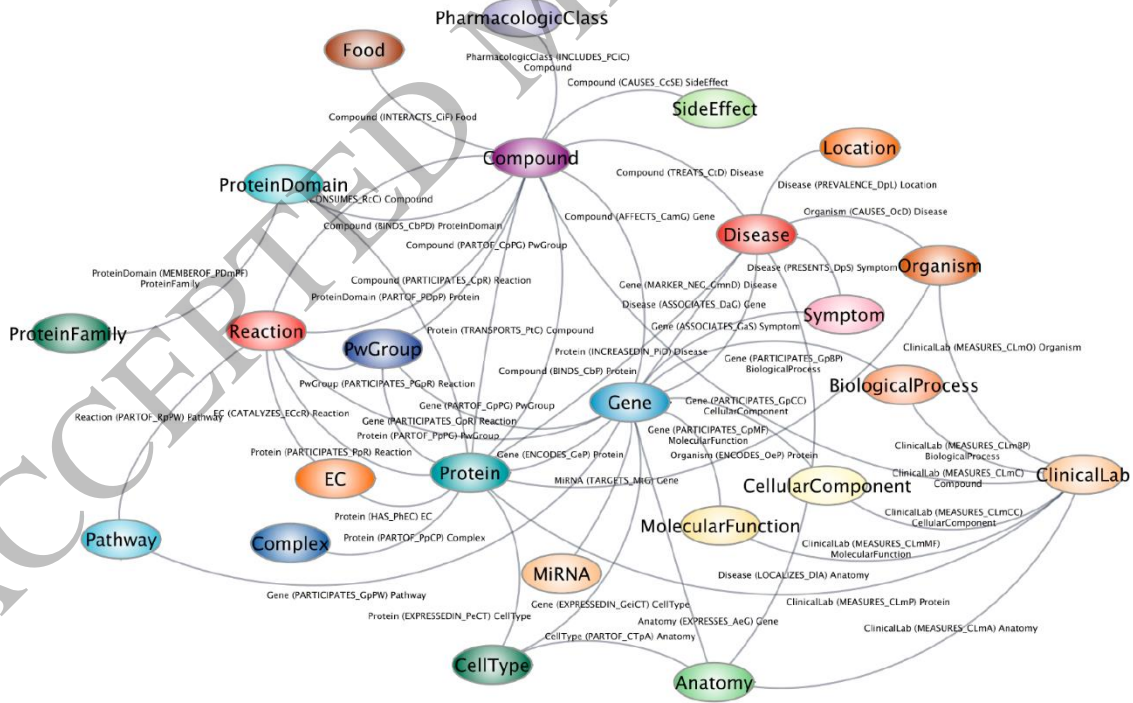
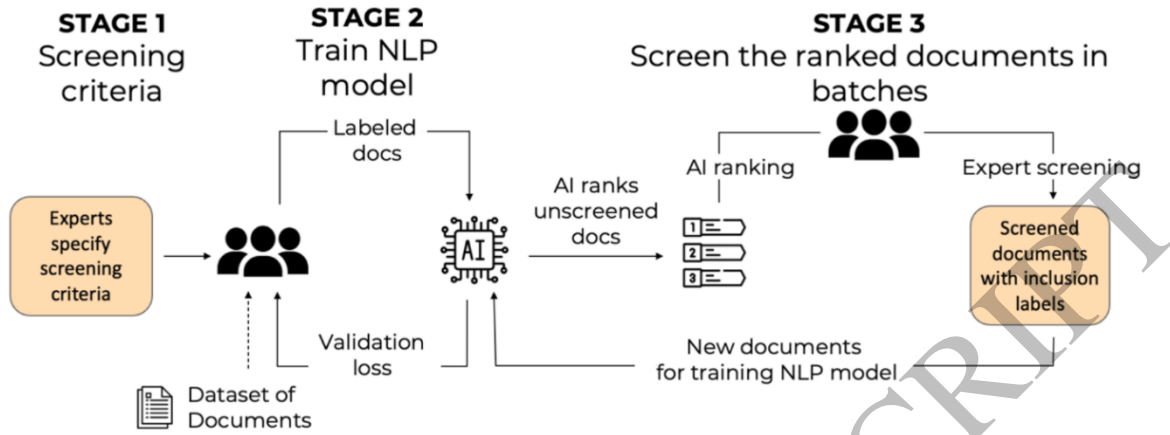
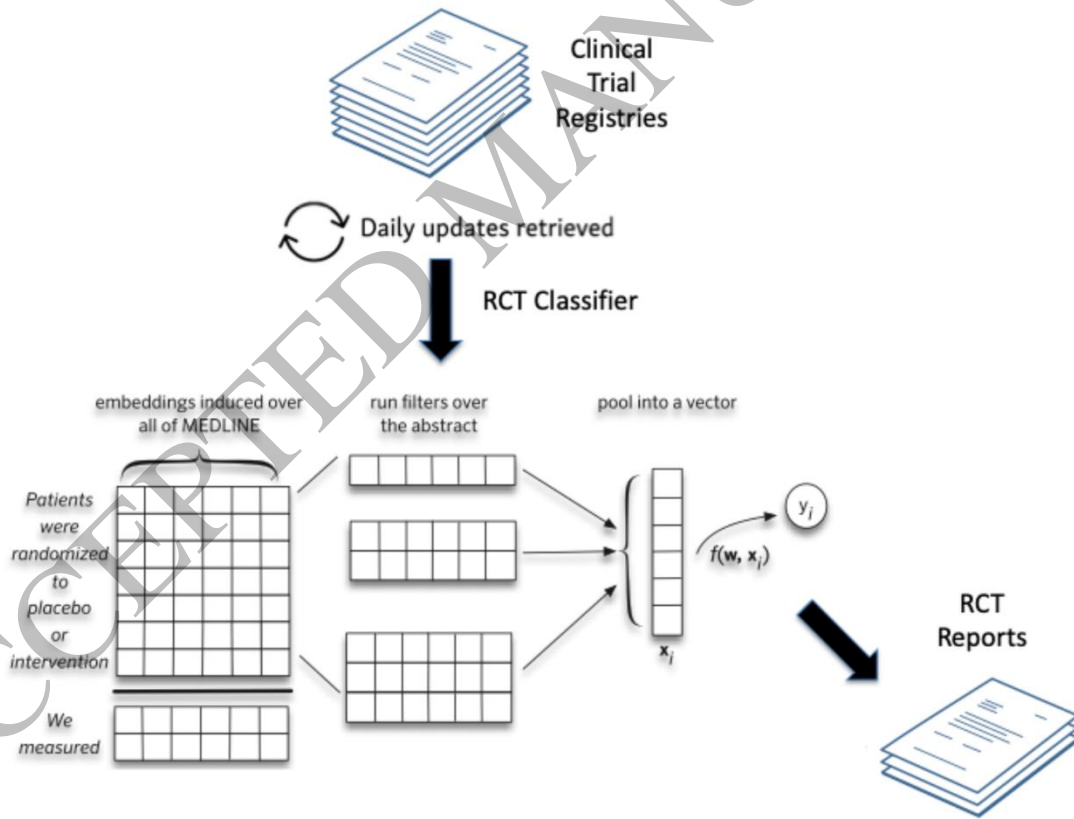


Figure 2
148x93 mm (x DPI)



1
2
3
4

Figure 3
150x56 mm (x DPI)



5
6
7

Figure 4
159x116 mm (x DPI)