
Evident: a Development Methodology and a Knowledge Base Topology for Data Mining, Machine Learning and General Knowledge Management

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Software has been developed for knowledge discovery, prediction and management
2 for over 30 years. However, there are still unresolved pain points when using
3 existing project development and artifact management methodologies. Historically,
4 there has been a lack of applicable methodologies. Further, methodologies that
5 have been applied, such as Agile, have several limitations including scientific
6 unfalsifiability that reduce their applicability. *Evident*, a development methodology
7 rooted in the philosophy of logical reasoning and *EKB*, a knowledge base topology,
8 are proposed. Many pain points in data mining, machine learning and general
9 knowledge management are alleviated conceptually. *Evident* can be extended
10 potentially to accelerate philosophical exploration, science discovery, education as
11 well as knowledge sharing & retention across the globe. *EKB* offers one solution
12 of storing information as knowledge, a granular level above data. Related topics in
13 computer history, software engineering, database, sensing hardware, philosophy,
14 and project & organization & military managements are also discussed.

15 1 Introduction

16 Necessity is the mother of invention, claimed Plato [1]. Deficient in rigorous scientific scrutinization
17 as the statement is, major methodology evolutions in software development did not emerge until the
18 emergence of major computer innovations and thereafter elevated effort orchestration needs.

19 In the 1940s, digital programmable electronic computers revolutionized scientific calculation done
20 previously with mechanical and analog computing machines [2]. Assembly (1947) [3, 4] and high
21 level (1953) [5, 6] programming languages rose to harness the unprecedented and ever-increasing
22 computing power. Eventually the term software was coined (1953) [7]. Two Software Development
23 Methodologies (SDMs) were proposed: 1. a project breakdown of sequential phases, the essence of
24 Waterfall (see Fig 1a), first presented no later than 1956 [8] and 2. iterative and incremental SDM,
25 the essence of Agile (see Fig 1b), first executed no later than 1957 [9].

26 In the 1960s, operating systems (1962) emerged to orchestrate multiple computation tasks[10]. This
27 signified the shift of computer development from single-task specialized machines for military and
28 academia to machines accessible to the general public. The shift was exemplified by The Mother of
29 All Demos (1968) which demonstrated many fundamental elements of personal computing for the
30 first time [11] and showed how software had evolved in both diversity and complexity. Meanwhile,
31 the first formal detailed diagram of the Waterfall methodology appeared in literature (See Fig 1a)
32 (1970) [12] and the name of Waterfall was ultimately coined (1976) [13]. Agile variants such as
33 evolutionary project management [14] and adaptive SDM [15] appeared in the early 1970s, although
34 no clear preference between Waterfall and Agile variants was found in the literature.

Table 1: Major Software Usage Evolutions and Methodology Developments

Period	Technology	Software Usage	Major Methodology Development
1940s	Programming Language	Scientific Calculation	First Waterfall variant presentation (1956) [8]; first Agile variant execution (1957)[9].
1960s	Operating System	Shifting to Applications	First detailed diagram of Waterfall idea(1970)[12], Waterfall name (1976)[13]; Agile variants: evolutionary project management [14] & adaptive SDM [15](early 1970s).
1980s	GUI & Internet	PC & Internet Applications	Waterfall standardized in military (1985)[20]; The Manifesto signed (2001)[24]. Agile significantly more popular than Waterfall.
Around 1990	Data Storage	Knowledge Discovery, Prediction and Management	NA

35

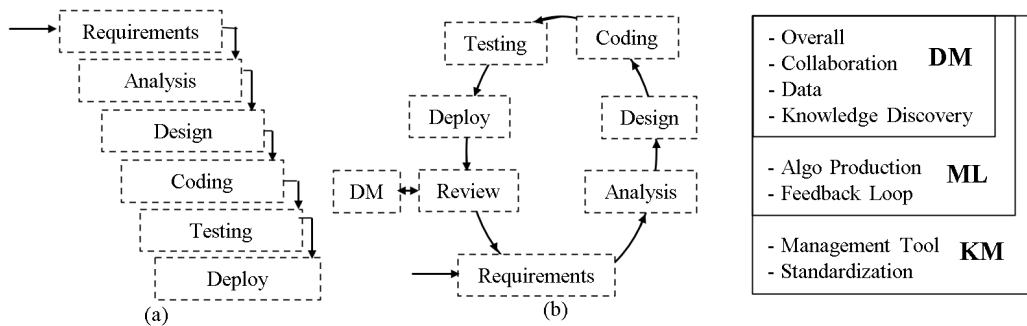


Fig 1 a) A typical Waterfall design diagram [12] b) An iterative, evolutionary and incremental design cycle commonly viewed as Agile[9] and how DM can improve each design cycle by discovering Knowledge about user needs during Review.

Fig 2. Relations among DM, ML and KM for their pain points

36

37 In the 1980s, personal computers entered households[16] followed by graphic user interface (GUI)
 38 (1983) [16]. The Internet Protocol Suite (TCP/IP) was standardized (1982)[17] and commercial
 39 Internet service providers emerged (1989) [18, 19] Unprecedented user-computer interactions and
 40 user-user communications created tremendous software needs, while Waterfall was still widely
 41 deployed in software development. United States Department of Defense issued a military standard
 42 describing Waterfall as the required military software development process (1985) [20]. However,
 43 software user needs grew so fast that, the heavy Waterfall SDM failed to deliver in pace. Consequently,
 44 a number of light weight SDMs were proposed and practiced (1990s) [4, 21, 22, 23]. Eventually The
 45 Manifesto for Agile Software Development (The Manifesto) [24] was signed by 17 practitioners of
 46 light-weight SDM (2001) and became the de facto SDM.

47 Around 1990, data storage capacities grew significantly and software usages in Data Mining (DM,
 48 defined as knowledge discovery from data) reached the tipping point. While data and software's
 49 storage manners differ between Von Neumann and Harvard architectures, data storage capacity growth
 50 empowered software to discover knowledge supported by scientific evidence (defined as Knowledge)
 51 that people never had access to. Corporations started to analyze customers' behavior and make
 52 business decisions based on Knowledge (1990s) [25]. The first DM methodology, Cross-Industry
 53 Standard Process for Data Mining (CRISP-DM) was conceived (1996) [26, 27]. However, CRISP-
 54 DM and its variants appear more of a theoretical framework, offer little meaningful or actionable
 55 guidance, and therefore have not gotten much traction. In addition, CRISP-DM is concerned only
 56 with DM, not Machine Learning (ML, defined as to deliver an algorithm (Algo) for Knowledge
 57 prediction) or Knowledge Management (KM) in general for science, medicine, military and so on.
 58 Due to the absence of alternative methodologies(see Table 1), Agile is still being offered up for DM,
 59 ML, and KM [28, 29] with questions being asked about its appropriateness [30, 31].

60 This paper discusses limitations in Agile as a scientific claim and why it may not address the current
 61 pain points of DM, ML and KM, which are later summarized. *Evident* along with *Evident Knowledge*
 62 *Base (EKB)* is proposed as a project development and artifact management methodology. *Evident's*
 63 potential in alleviating many current pain points is demonstrated conceptually. Unalleviated pain

64 points and future work to fulfill the potential are also discussed. Beyond software development,
 65 *Evident* is illustrated to be applicable in many aspects of society. *EKB* is demonstrated as one potential
 66 infrastructure to store information as Knowledge, a granular level above data.

67 2 Agile ambiguity and unfasifiability

68 Most people regard Agile as iterative, evolutionary and incremental software development [9] (see
 69 Fig 1b) and many claim to be Agile practitioners. However, Agile empirical evidence is mixed and
 70 hard to find [32, 30] while no measurable scientific evidence has been found at all. Although control
 71 experiment challenges or absence of quantitative project Agility measurements may explain no
 72 measurable scientific evidence, concerns remain with the ambiguity with which Agile’s approaches
 73 and scope are defined in The Manifesto.

74 2.1 Agile approaches are vaguely defined in The Manifesto

75 The Manifesto includes 4 values and 12 principles [24]. The goal is crystal clear: to rapidly deliver
 76 quality software that meets user needs, but not so much can be found for how to get there. Most of
 77 the Values and Principles appear to be goals but not approaches (see Appendix); some are concerned
 78 with approaches but vaguely defined; only four principles are actionable, which turn out to have no
 79 relevance in how to implement iterative, evolutionary or incremental development. Agile Alliance,
 80 co-founded by some original signers of The Manifesto, defines Agile as “an umbrella term for a set
 81 of frameworks and practices” from which Agile practitioners “figure out the right things to do given
 82 your particular context.” [33] Unfortunately, no actionable approaches are defined.

83 Therefore although there are numerous frameworks under the Agile umbrella [34, 35], it’s impossible
 84 to determine if a development practice is Agile and the claim of Agile practice becomes unfalsifiable.
 85 Because falsifiability is the standard evaluating scientific against non-scientific claims introduced by
 86 Karl Popper [36], Agile is not a scientific claim. Consequently, no observable scientific evidence can
 87 prove or disprove Agile, because technically no one can determine if a project is Agile or not in the
 88 first place. If an Agile rollout “fails”, Agile proponents can always argue that the Agile rollout was
 89 not implemented correctly.

90 Meanwhile Agile practitioners cannot determine if they practice Agile correctly either. Projects
 91 employing Agile frameworks such as Test Driven Development or Feature Driven Development
 92 may not even realize that the projects may not be adaptive to new user needs. People who are
 93 essentially practicing Waterfall may believe they are practicing Agile only because they implement
 94 each sequential Waterfall phase incrementally or simply use Scrum or Kanban.

95 In defense, some proponents claim Agile as a philosophy [37, 38]. Granted Agile’s goal may fit into
 96 Axiology, one of Philosophy’s four domains (the rest as Metaphysics, Epistemology and Logic) [39],
 97 concerned with what is good, it appears to be a common understanding and offers little value when
 98 Agile’s approaches are vaguely defined.

99 2.2 No scopes are defined in The Manifesto

100 “The right things to do given your particular context” by Agile Alliance [33] are expected to be found
 101 within Agile’s frameworks, otherwise Agile is not practiced right. With no scope defined, Agile
 102 seems to cover the scope of all softwares. However, some softwares have non-incremental needs or
 103 simply only one need, e.g. to solve one specific partial differential equation numerically. Their needs
 104 are either met or not at all. No iteration or evolutionary Agile design cycles exist.

105 Agile is also not applicable for Knowledge discovery tasks such as DM. In a typical Agile development
 106 cycle (Fig 1b), Review phase is to discover Knowledge about user needs, which can be done through
 107 DM. Therefore Agile should not be applicable to DM, one phase of its own design cycle.

Table 2: Pain Points of DM, ML and KM.

Pain Points for DM		
Overall	Struggles to deliver fast with technical debt	abc
	Project Progress not easily measurable	b
	Uncertainty in project timeline	b
	Activities not easily trackable or reproducible	bc

	Not scalable in terms of both collaborator number and project maintenance	ab
	Few general project design patterns	abc
	Anti-patterns not uncommon	abc
Collaboration	No methodologies to orchestrate team of size commonly seen in software development	b
	Few intuitive manners to divide task among team members	ab
	Deficient common awareness in needs for process improvement	a
	Tasks completed or ideas explored by team members cannot be easily found and reproduced causing duplicated work.	bc
Data	Data compromised in availability, accuracy and consistency during acquisition	x
	Data preprocessing process not standardized such as data labeling, object detection (e.g. identify object pixels in images) causing unstable data dependency, cascade correction and uncertainty in project progress	b
	Underutilized data may take unnecessary resource	b
	Data may be presented in different data type such as integer, float or string, causing unnecessary data dependency for Algo and experiments	b
Knowledge Discovery / ML Algo Research	Off-the-shelf models are available for DM automation. However, DM automation has not become a common practice.	abc
	Inefficient in-house model code implementation not uncommon	b
	No appropriate version control tools. Current version control tools such as git are designed to only keep the best version Algo/experiment available, while DM and ML need multiple versions available concurrently for reference	abc
	No easy solution to request flexible data storage, memory and computation capacity as needed. Hard drive, RAM, CPU and GPU are difficult to allocate even on the cloud.	x
	The use of other Algos' output as input results in correction cascades	b
	Algo is a sequential computation process different from typical software applications with a number of independent features. Difficult to assign one Algo development into multiple team members	x
	Algo user may have no clear understanding about the Algo and deploys it outside its scope.	b
	Multiple programming language smell	x
Pain Points for ML		
Algo Production	Significant efforts of research Algo migration into production	bc
	Even more significant efforts if production Algo is written in a different language than the language used in research, e.g. in embedded system	x
	For Algo analyzing sensor data such as cameras or bio-sensors, product grade data won't be available for Algo research until sensor hardware designs are complete. Algo becomes the product release bottleneck, resulting in either sub-optimal production Algo or delayed product release.	x
	Production data source is inconsistent with research data source	x
	Algo production is often done by team members, most likely software engineers, who did not produce the Algo, causing misuse	abc
	Algo needs to load data in real time during production but most often not in real time during research, causing unnecessary Algo code re-factoring	b
	Prototype Algo may be accidentally run in production causing damages	x
	No straightforward way to organize codes repository for research and production team members work in the same repository	abc
	Dead code path	b
Feedback Loop	Algo update workflow not straight forward	b
	Few clear pattern designs for monitoring Algo performance in production	b
	Actions based on unseen data predicted by Algo may alter observed data	x
Pain Points for KM		
Management Tool	Few tools or resource help people check if Knowledge formed is well supported by evidence, especially when evidence appears long after presumed Knowledge has been formed.	abc
	Knowledge dissemination among community has always been a challenge.	abc
	Knowledges formed by different organizations are not easy to combine	abc

	Knowledge formed within organizations is not easy to share and retain.	abc
Standardization	Knowledge has been recorded in sentences or articles. Few standardized ways to represent general Knowledge.	abc

a/b/c: Pain points that can be alleviated by *Evident*'s character a, b or c. x: Pain points that cannot be alleviated by *Evident*.

108 3 Pain points for DM, ML and KM

109 Owing to the absence of applicable methodologies, pain points have been continuously reported for
 110 DM, ML and KM [40, 41, 42, 43] (see Table 2) in the current big data era with explosive growth in
 111 data volume, variety and velocity.

112 DM & ML's Algo typically comprises a data computation flow (defined as a Model, supervised
 113 or unsupervised), such as logistic regression, and its configuration, such as logistic regression
 114 coefficients. DM employs off-the-shelf or in-house Models to discover Knowledge from data. ML
 115 compares Knowledges discovered by DM by candidate models and deploys the one that performs best
 116 with its configuration, as the production Algo, for Knowledge prediction in production. Therefore,
 117 DM's pain points still apply to ML. Meanwhile because DM and ML are special forms of KM, their
 118 pain points are also applicable for KM (see Fig 2).

119 Generally speaking, DM has not been regarded highly collaborative and scalable activities to deliver
 120 high throughput Knowledge. It's rare to see hundreds of contributors in a DM project, unlike for
 121 example some complicated open source software projects that, deliver promptly, efficiently and
 122 continuously for years or even decades [44, 45]. It is also rare of mass Knowledge production
 123 in a organized and standardized manner with high production yield for a unit period, commonly
 124 seen in consumer products such as automobiles or toothpastes. DM often struggles to deliver
 125 Knowledge rapidly with technical debts in reproducibility, measurability, trackability. DM needs
 126 to not only handle artifacts of different modalities such as documents, codes and data, but also
 127 address computation and data storage resource requests potentially across multiple platforms. Raw or
 128 preprocessed data can be compromised in availability, accuracy and consistency. Data dependency
 129 and entangled models often cause cascaded correction and uncertainty in project planning. In addition,
 130 routine tasks such as quarterly or annual finance analysis mostly have not be automated. Tools and
 131 project management methodologies are highly in demand to fulfill DM's potential and deliver values.

132 In ML, the team members, usually software engineers or product managers, that deploy an Algo in
 133 production to predict future data may not have produced the Algo and may misuse it. The Algo codes
 134 are often refactored sometimes in different programming languages, operating systems or even in
 135 fixed point instead of float point. If the Algo is to be deployed on a data acquisition product such
 136 as cameras or bio-sensors, no product grade data are available for Algo research until the sensor
 137 hardware design is finalized to enable data collection. Algo research therefore becomes the bottleneck
 138 for product release. Frequently sub-optimal Algo is deployed to meet the deadline or projects become
 139 delayed. What's more, production data may come from different sources compared to research data
 140 potentially caused by, e.g. sensor upgrade or downgrade, resulting in the under performance of
 141 production Algo. After Algo deployment, no straightforward way exists to monitor Algo performance
 142 or thereafter update Algo. Future Knowledge predicated by the deployed Algo may encourage
 143 Algo users to alter decisions, which leads to the formation of future data with unanticipated hidden
 144 feedback loops. In addition, prototype Algo or dead codes may be accidentally run in production
 145 potentially causing catastrophic consequences.

146 Although scientific methods have guided people to discover Knowledge and improve practices such
 147 as evidence-based medicine [46] and experiment based military development [47], people still form
 148 Knowledge that lacks in supporting evidence [48]. One possible reason is the unavailability of tools
 149 or resources to check if the Knowledge formed is well supported by evidence, especially when there
 150 is a significant time gap between the Knowledge formed and the appearance of supporting evidence,
 151 e.g. for long-term investments or corporation strategies. Knowledge dissemination and retention
 152 are also huge challenges among the community and organization [49]. Furthermore, Knowledge is
 153 mostly recorded in the form of articles. However, because articles writing has not been and probably
 154 will never be standardized, Knowledge has not been able to be represented in a standardized manner
 155 for definition, reference and storage. The same Knowledge recorded in different sentences or even
 156 languages may be interpreted differently.

157 **4 Evident: a project development and artifact management methodology**

158 **4.1 Definition and scope**

159 *Evident* is a methodology of project, including but not limited to software, development and artifact
 160 management for DM, ML and KM, characterized by

- 161 a. project development mimicking a continuous process of logical reasoning in philosophy;
- 162 b. project activities or artifacts are broken into containers of Observations, Hypotheses and
 163 Tests (collectively defined as Containers);
- 164 c. directional association constructions towards and only towards Test Containers to represent
 165 Knowledge.

166 Observation is a collection of facts. A Hypothesis is Knowledge to be formed out of Observation. A
 167 Test is a Hypothesis evaluation process using Observation to prove or disprove the Hypothesis with
 168 or without confidence levels. Containers indicate Observations, Hypotheses and Tests can only be
 169 added or removed as a block.

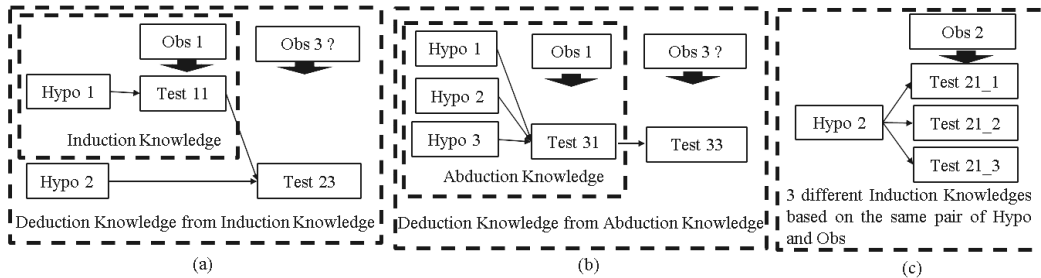


Fig 3 Induction, abduction and deduction Knowledge represented in *Evident* (Hypo: Hypothesis; Obs: Observation)

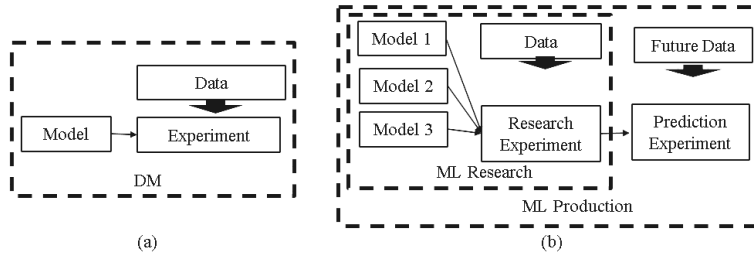


Fig 4 DM, ML processes represented in *Evident*. a) DM mimic Knowledge induction in Fig 3a; b)
 ML Research and Production mimic Knowledge abduction and Knowledge deduction in Fig 3b.

171 **4.2 Knowledge representation**

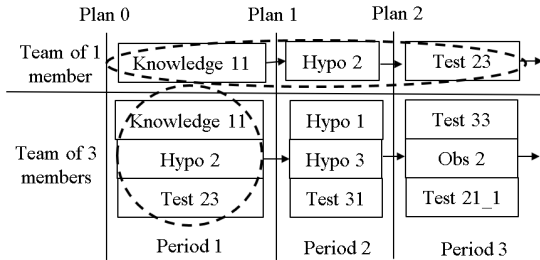
172 Loosely speaking, a Test associated with a Hypothesis and an Observation represents induction
 173 Knowledge (see Fig 3a); a Test associated with a Hypothesis set and an Observation represents
 174 abduction Knowledge (see Fig 3b), a Hypothesis associated Test that is also associated with an
 175 induction or abduction Knowledge Test represents deduction Knowledge or prediction (see Fig 3a&b).
 176 Deduction Knowledge becomes induction Knowledge once Observation proving or disproving
 177 deduction Knowledge is associated with Test, while stay deducted Knowledge if the associated
 178 Observation overlooks (fails to either prove or disprove), the deduction Knowledge. Multiple
 179 Tests associated with the same pair of Hypothesis(es) and Observation represent multiple different
 180 Knowledges based on different evaluation metrics (e.g. profit maximization or cost minimization) or
 181 Observation usage strategies (e.g. cross-validation grouping) (see Fig 3c).

182 DM may be regarded as Knowledge induction with data as Observation, model as Hypothesis to
 183 be evaluated and data analysis experiment as model test on data to form induction Knowledge with
 184 statistical confidence (see Fig 4a). Similarly, ML is Knowledge abduction. An example is a data
 185 analysis experiment that picks the best off-the-shelf or in-house model that best explains the data
 186 to form the Algo for production (see Fig 4b). Experiments employing different cost functions (e.g.

187 RMSE, AUC or correlation coefficients), statistical confidence levels or data allocation strategies for
 188 training and testing may result in different Knowledges or Algos.

189 4.3 Project development

190 *Evident* project developments are intuitively broken down into two granular levels: Knowledges
 191 and Containers. Mimicking logical reasoning in philosophy, each project period develops a batch
 192 of independent Knowledges or Containers assigned to teams of various sizes to maximize unit time
 193 throughput (see Fig 5). The next batches of Knowledge or Containers can be adaptively planned
 194 after period reviews or retrieved from backlogs. *Evident* is compatible with Kanban, Scrum or other
 195 development tools or frameworks for project planning and development of a single Knowledge or
 196 Container. Any tools or frameworks that do not compromise the project breakdown into Knowledge
 197 and Containers are applicable.



198 Fig. 5 One exemplary *Evident* project development process to produce Knowledges in Fig 3 a&b&c. Circled tasks may take 1 team member 3 periods, or take a team of 3 members 1 period. Knowledge 11 represents Obs 1 & Hypo 1 & Test 11 in Fig 3a.

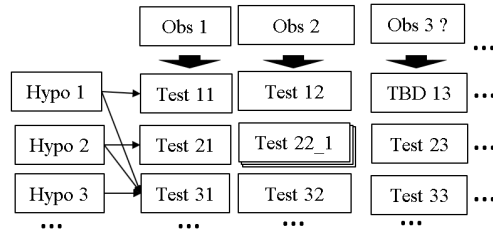


Fig. 6 An exemplary EKB composed of Hypo, Obs and Test containers that can store induction, abduction and deduction Knowledges in Fig 3 a&b&c and project development artifacts in Fig 5 continuously.

199 4.4 Artifact management: EKB

200 *Evident* artifact management may build on a topology of a relational Knowledge base, named
 201 *EKB*, composed of *Evident* Containers and directional associations (see Fig 6). Knowledges can
 202 be reproduced by Containers stored in *EKB*. Oversimplified as a table, *EKB* columns represent
 203 Observations; rows represent Hypotheses; values represent Tests or Test to be done (TBD). A Test
 204 can only be associated with one Observation, even if the associated Observation overlooks the
 205 Hypothesis(es) associated with the Test. Any new Observation proving, disproving or overlooking
 206 the same Hypothesis(es) occupies a column in *EKB*.

207 4.4.1 EKB stores containers and Knowledges continuously

208 When new Hypotheses, Observations or deduction Knowledge Tests are developed, new rows or
 209 columns of TBDs are inserted. A Test associated with Hypothesis(es) and a Observation can be
 210 stored in the designated row and column to represent different Knowledges.

211 An induction Knowledge Test is placed in the row of its associated Hypothesis and the column of its
 212 associated Observation (see Fig 3a&c & Fig 6). An abduction Knowledge Test is placed in the row
 213 of the Hypothesis best explaining the Observation (see Fig 3b & Fig 6). A deduction Knowledge
 214 is placed in the row of its associated Hypothesis and the column of the pending Observation (see
 215 Fig 3a&b & Fig 6). Once an Observation becomes available proving or disproving the Hypothesis,
 216 a deduction Knowledge becomes an induction Knowledge. Multiple Tests representing different
 217 Knowledges can be placed in the same slot (see Fig 3c & Fig 6).

218 4.4.2 EKB supports relational database operations of Permutation And Join

219 *EKB* is similar to a relational database [50] with Observations as columns, Hypotheses as rows and
 220 Tests as values, but with potential associations among values for deduction Knowledges. Relational
 221 database operations independent of values associations, such as Permutation (switching rows and
 222 columns) and Join (merging *EKBs*) can be implemented without compromise in *EKB*; operations
 223 dependent on values associations such as Restriction (select rows), Projection (select columns) and
 224 Compositions (merge selected columns&rows from multiple *EKBs*) can only be implemented for
 225 *EKBs* storing only induction or abduction Knowledge and have no associations among Tests.

226 *EKBs* are highly flexible for team collaboration and maintenance. Different team members working
227 on different Containers can share one *EKB* as the common work space to improve efficiency. Multiple
228 *EKBs* can be joined together without information loss so that Knowledges produced by different
229 teams or team members can be accumulated into one *EKB*. For *EKBs* storing only induction and
230 abduction Knowledge, all relational database operations are applicable, so that team members can
231 compose their own *EKBs* without keeping a potentially large team *EKB* on the local machines.

232 **5 Advantages and pain points alleviated**

233 Inspired by the philosophy of logical reasoning, *Evident* is intuitive to understand and follow. Project
234 activity and artifact containerization supports incremental as well as adaptive project planning and
235 artifact pattern abstraction. Disentangling Hypotheses and Observations reduces unnecessary de-
236 pendency, cascade correction and uncertainty in project planning. Knowledge representation in
237 associations among artifacts can not only track Knowledge development, but also Knowledge devel-
238 opment status (prove, disproved or overlooked), which improves project measurability, trackability
239 and reproducibility. Overall *Evident* may help applicable projects deliver fast and at scale with many
240 pain points alleviated in DM, ML and KM (see Table 2).

241 **5.1 DM**

242 Containerized Data and Models in *Evident* prevents unstable data dependency, model entanglement
243 and cascade correction. Dead data and codes can be easily identified and removed. Standardized
244 Models and Experiments encourage reuse of computationally efficient containers, support automatic
245 DM. Different experiments may use different optimization target function on the same Model and
246 Data to deliver different Knowledges for different users, e.g. Marketing vs Engineering managers.

247 Containerization is an alternative to the state of the art artifact version control, such as git, which
248 keeps only the one version of the code or data in the workspace with historic versions saved as
249 commits. *Evident* keeps all applicable versions available in the workspace for easy access. This may
250 appear to use more storage space. However current version control tools all save version commits as
251 snapshots [51], demanding comparable storage space of *Evident* if a *Evident* equivalent number of
252 versions are stored.

253 *Evident* granulates project activities into independent standardized Knowledge and Container levels,
254 supports adaptive development, facilitates project planning among collaborators in teams of various
255 sizes and reduces planning overhead. Artifacts are continuously stored in *EKB*, making project
256 development measurable, trackable, reproducible and scalable. Meanwhile once Containers are
257 produced, Knowledge or documentation reports can be generated automatically instead of manually.
258 *Evident* accelerates DM delivery in both short-term and long-term.

259 **5.2 ML**

260 A deployed Algo can be evaluated easily by re-applying the original Research Experiment on the
261 production data. Different users involved in the deployment can understand the Algo's scope and
262 origins easily by examining the research Experiment (see Fig 4b). Research Experiments can load
263 data in real time as Production Experiment, so that both Experiments can inherit the same design
264 patterns with statistical analysis and evaluation metrics. The Production Experiment can report and
265 examine the prediction performance at regular time intervals to detect production data pattern drift
266 for either model reconfiguration or model replacement. Once a model with its configuration is retired
267 from production, the production data is containerized and associated with the Production Experiment,
268 transforming the Production Experiment into Research Experiment and a deduction Knowledge for
269 prediction into an induction Knowledge that is also preserved in *EKB*.

270 Because both research and production can operate on the same *EKB*, research and production team
271 members can share the same workspace the way software engineers work on the same code repository,
272 facilitating the model migration from research to production and efficient team collaborations.

273 5.3 KM

274 Knowledge formatting into design patterns of Containers provides a meaningful progress towards
275 Knowledge standardization for improved definition, reference and storage compared to state of art
276 sentences or articles. *EKB* with standardized Container templates may offer potential tools for people
277 to examine the Hypotheses formed against evidence or Observations, facilitating evidence-based
278 decision making and Knowledge development. *EKB* can not only facilitate Knowledge dissemination,
279 accumulation and retention, but also label the development status of each Hypothesis as proved,
280 disproved or overlooked, a desirable design pattern for projects and Knowledge Management.

281 6 Discussions

282 6.1 Significance

283 *Evident* may advance many society domains such as software, philosophy, science, business as well
284 as Knowledge sharing and retention across the globe, thanks to its applicability to general KM.

285 *EKB* may make no smaller impacts than relational data base [50], the invention of which created a
286 data base industry, as one solution to store information as Knowledge, a granular level above data.

287 6.2 Work to do

288 Work needs to be done regarding *Evident* ergodicity over logical reasoning in philosophy. If proved,
289 *Evident* can support all logical reasoning in philosophy. No evidence has existed to prove or disprove
290 the claim. *Evident* ergodicity is overlooked, stated in *Evident* language, especially considering logical
291 reasoning in philosophy may evolve.

292 Control studies need to be done to show *Evident* can truly provide value. No tools tailored to support
293 *Evident* project development planning and *EKB* are available, although some existing tools are
294 applicable for use. Particularly the tools that allow unexpected alteration proof, easy access and
295 visualization of Containers are in demand. More detailed discussions need to be done about how
296 *Evident* help applicable projects with examples.

297 6.3 Pain points not alleviated

298 *Evident* offers no detailed guidance in development below Container level. For example, a Model or
299 computation flow cannot be broken down further into smaller modules by *Evident* for incremental and
300 adaptive development. Multiple languages smells and accidents running prototype Algo in production
301 cannot be avoided by *Evident* either. In addition, *Evident* cannot control future observation alteration
302 caused by decisions made by people based on *Evident* produced Knowledge.

303 *Evident* can only manage project artifacts of data or codes but not sensors or hardwares. Pain points
304 caused in data acquisition such as availability, inaccuracy and inconsistency are out of *Evident*'s
305 scope. *Evident* is incapable of improving computation hardware resources allocation either.

306 6.4 More words about Agile

307 Due to Agile's ambiguity and unfalsifiability as a scientific claim, it might be a better practice to drop
308 the term Agile and instead quote each framework currently under Agile on its own. Frameworks
309 such as iterative and evolutionary development as well as Kanban are valuable although need to be
310 employed discretionally. Practitioners should have better understood what exactly they were doing
311 without being fuzzed by the buzzword Agile.

312 7 Conclusions

313 The paper proposes *Evident* as a project development and artifact management methodology for DM,
314 ML and KM as well as *EKB* as a Knowledge base topology. *Evident* and *EKB* have been shown
315 of great value to alleviate many unresolved pain points. *Evident* has the potential to facilitate the
316 advancement of many aspects of society due to its utility in general Knowledge management. *EKB*
317 may serve as the infrastructure for storing information as Knowledge, a granular level above data.

318 **A Appendix**

319 **A.1 The following among the four Values and twelve Principles of The Manifesto [24] appear**
320 **to be goals:**

- 321 Value 4: Responding to change over following a plan;
- 322 Principle 1: Customer satisfaction by early and continuous delivery of valuable software.
- 323 Principle 2: Welcome changing requirements, even in late development.
- 324 Principle 3: Deliver working software frequently (weeks rather than months)
- 325 Principle 7: Working software is the primary measure of progress
- 326 Principle 8: Sustainable development, able to maintain a constant pace
- 327 Principle 9: Continuous attention to technical excellence and good design
- 328 Principle 10 : Simplicity—the art of maximizing the amount of work not done—is essential

329 **A.2 The following in The Manifesto appear to be approaches but vaguely defined:**

- 330 Value 1: Individuals and interactions over processes and tools
- 331 Value 2: Working software over comprehensive documentation
- 332 Value 3: Customer collaboration over contract negotiation
- 333 Principle 5: Projects are built around motivated individuals, who should be trusted
- 334 Principle 11: Best architectures, requirements, and designs emerge from self-organizing teams

335 **A.3 The following in The Manifesto appear to be actionable approaches but irrelevant of**
336 **iterative, evolutionary or incremental development regarded as Agile by most people [9]:**

- 337 Principle 4: Close, daily cooperation between business people and developers
- 338 Principle 6: Face-to-face conversation is the best form of communication (co-location)
- 339 Principle 12: Regularly, the team reflects on how to become more effective, and adjusts accordingly

340 **References**

- 341 [1] B. Jowett and L. Campbell, *Plato's Republic : the Greek text.* Oxford : At the Clarendon
342 Press, 1894.
- 343 [2] J. J. P. Eckert and J. W. Mauchly, "Electronic numerical integrator and computer," New York,
344 NY, 1947.
- 345 [3] A. D. Booth and K. H. Britten, "General considerations in the design of an all purpose electronic
346 digital computer," Tech. Rep., 1947.
- 347 [4] M. Campbell-Kelly, "The development of computer programming in Britain (1945 to 1955),"
348 *Annals of the History of Computing*, vol. 4, no. 2, pp. 121–139, 1982.
- 349 [5] P. Bentley, *Digitized: The science of computers and how it shapes our world.* New York:
350 Oxford University Press, 2012.
- 351 [6] D. E. Knuth and L. T. Pardo, "Early development of programming languages," *A History of*
352 *Computing in the Twentieth Century*, 1980.
- 353 [7] R. R. Carhart, *A Survey of the Current Status of the Electronic Reliability Problem.* Santa
354 Monica, CA: RAND Corporation, 1953.
- 355 [8] *Symposium on advanced programming methods for digital computers : Washington, D.C., June*
356 *28, 29, 1956.* Office of Naval Research, Dept. of the Navy, 1956.

- 357 [9] C. Larman and V. R. Basili, "Iterative and incremental development: A brief history," *Computer*,
358 vol. 36, pp. 47–56, 2003.
- 359 [10] S. H. Lavington, *A History of Manchester Computers*. British Computer Society, 1998.
- 360 [11] T. Bardini, *Bootstrapping: Douglas Engelbart, Coevolution, and the Origins of Personal*
361 *Computing*. Stanford CA: Stanford University Press, 2000.
- 362 [12] W. W. Royce, "Managing the development of large software systems," in *Technical Papers of*
363 *Western Electronic Show and Convention*, 1970.
- 364 [13] T. E. Bell and T. A. Thayer, "Software requirements: Are they really a problem?" in *Proceedings*
365 *of the 2nd International Conference on Software Engineering*. Washington, DC, USA: IEEE
366 Computer Society Press, 1976.
- 367 [14] T. Gilb, *Software Metrics*. Winthrop Publishers, 1976.
- 368 [15] E. A. Edmonds, "A process for the development of software for nontechnical users as an adaptive
369 system," *General Systems*, vol. 19, p. 215–18, 1974.
- 370 [16] W. Isaacson, *The Innovators: How a Group of Hackers, Geniuses, and Geeks Created the*
371 *Digital Revolution*. SIMON and SCHUSTER, 2014.
- 372 [17] J. L. Pelkey, "The history of computer communications," April 2022. [Online]. Available:
373 <https://historyofcomputercommunications.info/>
- 374 [18] R. Clarke, "Origins and nature of the internet in australia," *Emergence: Complexity and*
375 *Organization*, vol. 4, pp. 1990–1994, 2004.
- 376 [19] R. H. Zakon, "Hobbes' internet timeline 25," April 2022. [Online]. Available:
377 <https://www.zakon.org/robert/internet/timeline/>
- 378 [20] "Military standard: Defense system software development by department of defense," Jun 1985.
- 379 [21] J. M. Kerr and R. Hunter, *Inside RAD: How to Build a Fully Functional System in 90 Days or*
380 *Less*. McGraw-Hill, 1993.
- 381 [22] R. Nagel and R. Dove, *21st Century Manufacturing Enterprise Strategy: An Industry-Led View*.
382 Diane Pub Co, 1991.
- 383 [23] A. Presley, J. Mills, and D. Liles, "Agile aerospace manufacturing," 1995.
- 384 [24] M. Beedle, A. van Bennekum, and A. Cockburn, "Manifesto for agile software development,"
385 April 2001. [Online]. Available: <https://agilemanifesto.org/>
- 386 [25] D. J. Power, "A brief history of decision support systems," April 2022. [Online]. Available:
387 <http://DSSResources.COM/history/dsshistory.html>
- 388 [26] R. Wirth and J. Hipp, "Crisp-dm: Towards a standard process model for data mining," *Proceed-*
389 *ings of the 4th International Conference on the Practical Applications of Knowledge Discovery*
390 *and Data Mining*, 2000.
- 391 [27] IBM, "Crisp-dm help overview," April 2022. [Online]. Available: [https://www.ibm.com/docs/](https://www.ibm.com/docs/en/spss-modeler/SaaS?topic=dm-crisp-help-overview)
392 [en/spss-modeler/SaaS?topic=dm-crisp-help-overview](https://www.ibm.com/docs/en/spss-modeler/SaaS?topic=dm-crisp-help-overview)
- 393 [28] Microsoft, "Agile development of data science projects," April 2022. [Online]. Available:
394 <https://docs.microsoft.com/en-us/azure/architecture/data-science-process/agile-development>
- 395 [29] "Agile data science by data science process alliance," April 2022. [Online]. Available:
396 <https://www.datascience-pm.com/agile-data-science/>
- 397 [30] I. Lee, "6 reasons why i think agile data science does not
398 work," April 2022. [Online]. Available: [https://towardsdatascience.com/](https://towardsdatascience.com/6-reasons-why-i-think-agile-data-science-does-not-work-ee4dd680bb59)
399 [6-reasons-why-i-think-agile-data-science-does-not-work-ee4dd680bb59](https://towardsdatascience.com/6-reasons-why-i-think-agile-data-science-does-not-work-ee4dd680bb59)

- 400 [31] E. Yan, “Data science and agile (what works, and what doesn’t,” April 2022. [Online]. Available:
401 <https://eugeneyan.com/writing/data-science-and-agile-what-works-and-what-doesnt/>
- 402 [32] T. Dyba and T. Dingsøy, “Empirical studies of agile software development: A systematic
403 review,” *Information and Software Technology*, vol. 50, no. 9, pp. 833–859, 2008.
- 404 [33] “Agile 101 by agile alliance,” April 2022. [Online]. Available: [https://www.agilealliance.org/
405 agile101/](https://www.agilealliance.org/agile101/)
- 406 [34] P. Abrahamsson, O. Salo, J. Ronkainen, and J. Warsta, “Agile software development methods:
407 Review and analysis,” 2017.
- 408 [35] A. L. Fruhling and A. E. Tarrell, “Best practices for implementing agile methods: A guide for
409 department of defense software developers,” *Information Systems and Quantitative Analysis
410 Faculty Publications*, vol. 27, 2007.
- 411 [36] K. Popper, *The Logic of Scientific Discovery*. London and New York: Routledge, 2002.
- 412 [37] S. Braams, “The software development landscape: A rationalization of agile software
413 development as a strategy in the face of organizational complexity,” April 2022. [Online].
414 Available: <http://essay.utwente.nl/80784/>
- 415 [38] S. Forum, “Agile - methodology or framework or philosophy,”
416 April 2022. [Online]. Available: [https://www.scrum.org/forum/scrums-forum/6117/
417 agile-methodology-or-framework-or-philosophy](https://www.scrum.org/forum/scrums-forum/6117/agile-methodology-or-framework-or-philosophy)
- 418 [39] T. Schick, *Doing Philosophy: An Introduction Through Thought Experiments*. Mcgraw-Hill,
419 2009.
- 420 [40] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J.-F.
421 Crespo, and D. Dennison, “Hidden technical debt in machine learning systems,” in *Advances in
422 Neural Information Processing Systems*, vol. 28. Curran Associates, Inc., 2015.
- 423 [41] H. H. Nithya Sambasivan, Shivani Kapania, ““everyone wants to do the model work, not the
424 data work”: Data cascades in high-stakes ai,” 2021.
- 425 [42] T. J. Sebastian Schelter, Felix Biessmann, “On challenges in machine learning model manage-
426 ment,” *IEEE Data Eng. Bull.*, 2018.
- 427 [43] F. Kumeno, “Software engineering challenges for machine learning applications: A literature
428 review,” *Intelligent Decision Technologies*, vol. 13, no. 4, pp., vol. 13, no. 4, pp. 463–476, 2019.
- 429 [44] “Python,” April 2022. [Online]. Available: <https://www.python.org/>
- 430 [45] “Ubuntu,” April 2022. [Online]. Available: <https://ubuntu.com/>
- 431 [46] D. L. Sackett, “Evidence-based medicine,” *Seminars in Perinatology*, vol. 21, no. 1, pp. 3–5,
432 1997.
- 433 [47] S. Spear and T. Hone, “Succeeding in periods of change,” *Proceedings U.S. Naval Institute*,
434 March 2022.
- 435 [48] LessWrong, “Welcome to lesswrong!” April 2022. [Online]. Available: [https://
436 www.lesswrong.com/posts/bJ2haLkcGeLtTWaD5/welcome-to-lesswrong](https://www.lesswrong.com/posts/bJ2haLkcGeLtTWaD5/welcome-to-lesswrong)
- 437 [49] E. M. Rogers, *Diffusion of innovations*, 5th ed. New York, NY: Free Press, 2003.
- 438 [50] E. F. Codd, “A relational model of data for large shared data banks,” *Commun. ACM*, vol. 13,
439 no. 6, p. 377–387, jun 1970.
- 440 [51] Git, “Git internals - git objects,” April 2022. [Online]. Available: [https://book.git-scm.com/
441 book/en/v2/Git-Internals-Git-Objects](https://book.git-scm.com/book/en/v2/Git-Internals-Git-Objects)

442 Checklist

443 The checklist follows the references. Please read the checklist guidelines carefully for information on
444 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
445 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
446 the appropriate section of your paper or providing a brief inline description. For example:

- 447 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 448 • Did you include the license to the code and datasets? **[No]** The code and the data are
449 proprietary.
- 450 • Did you include the license to the code and datasets? **[N/A]**

451 Please do not modify the questions and only use the provided macros for your answers. Note that the
452 Checklist section does not count towards the page limit. In your paper, please delete this instructions
453 block and only keep the Checklist section heading above along with the questions/answers below.

- 454 1. For all authors...
 - 455 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
456 contributions and scope? **[Yes]**
 - 457 (b) Did you describe the limitations of your work? **[Yes]**
 - 458 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]**
 - 459 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
460 them? **[Yes]**
- 461 2. If you are including theoretical results...
 - 462 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 463 (b) Did you include complete proofs of all theoretical results? **[Yes]**
- 464 3. If you ran experiments...
 - 465 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
466 mental results (either in the supplemental material or as a URL)? **[N/A]**
 - 467 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
468 were chosen)? **[N/A]**
 - 469 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
470 ments multiple times)? **[N/A]**
 - 471 (d) Did you include the total amount of compute and the type of resources used (e.g., type
472 of GPUs, internal cluster, or cloud provider)? **[N/A]**
- 473 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - 474 (a) If your work uses existing assets, did you cite the creators? **[N/A]**
 - 475 (b) Did you mention the license of the assets? **[N/A]**
 - 476 (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
477
 - 478 (d) Did you discuss whether and how consent was obtained from people whose data you're
479 using/curating? **[N/A]**
 - 480 (e) Did you discuss whether the data you are using/curating contains personally identifiable
481 information or offensive content? **[N/A]**
- 482 5. If you used crowdsourcing or conducted research with human subjects...
 - 483 (a) Did you include the full text of instructions given to participants and screenshots, if
484 applicable? **[N/A]**
 - 485 (b) Did you describe any potential participant risks, with links to Institutional Review
486 Board (IRB) approvals, if applicable? **[N/A]**
 - 487 (c) Did you include the estimated hourly wage paid to participants and the total amount
488 spent on participant compensation? **[N/A]**