

# 1 Mimicking Clinical Trials with Synthetic Acute Myeloid Leukemia 2 Patients Using Generative Artificial Intelligence

3 Jan-Niklas Eckardt,<sup>1,2</sup> Waldemar Hahn,<sup>3,4</sup> Christoph Röllig,<sup>1</sup> Sebastian Stasik,<sup>1</sup> Uwe Platzbecker,<sup>5</sup>  
4 Carsten Müller-Tidow,<sup>6</sup> Hubert Serve,<sup>7</sup> Claudia D. Baldus,<sup>8</sup> Christoph Schliemann,<sup>9</sup> Kerstin Schäfer-  
5 Eckart,<sup>10</sup> Maher Hanoun,<sup>11</sup> Martin Kaufmann,<sup>12</sup> Andreas Burchert,<sup>13</sup> Christian Thiede,<sup>1</sup> Johannes  
6 Schetelig,<sup>1</sup> Martin Sedlmayr,<sup>4</sup> Martin Bornhäuser,<sup>1,14,15</sup> Markus Wolfien,<sup>3,4</sup> and Jan Moritz Middeke<sup>1,2</sup>

7 <sup>1</sup> Department of Internal Medicine I, University Hospital Carl Gustav Carus, Technical University Dresden, Dresden,  
8 Germany

9 <sup>2</sup> Else Kröner Fresenius Center for Digital Health, Technical University Dresden, Dresden, Germany

10 <sup>3</sup> Center for Scalable Data Analytics and Artificial Intelligence (ScaDS.AI) Dresden/Leipzig, Germany

11 <sup>4</sup> Institute for Medical Informatics and Biometry, Technical University Dresden, Dresden, Germany

12 <sup>5</sup> Medical Clinic and Policlinic I Hematology and Cell Therapy. University Hospital, Leipzig, Germany

13 <sup>6</sup> Department of Medicine V, University Hospital Heidelberg, Heidelberg, Germany

14 <sup>7</sup> Department of Medicine 2, Hematology and Oncology, Goethe University Frankfurt, Frankfurt, Germany

15 <sup>8</sup> Department of Hematology and Oncology, University Hospital Schleswig Holstein, Kiel, Germany

16 <sup>9</sup> Department of Medicine A, University Hospital Münster, Münster, Germany

17 <sup>10</sup> Department of Internal Medicine V, Paracelsus Medizinische Privatuniversität and University Hospital Nürnberg,  
18 Nürnberg, Germany

19 <sup>11</sup> Department of Hematology, University Hospital Essen, Essen, Germany

20 <sup>12</sup> Department of Hematology, Oncology and Palliative Care, Robert-Bosch-Hospital, Stuttgart, Germany

21 <sup>13</sup> Department of Hematology, Oncology and Immunology, Philipps-University-Marburg, Marburg, Germany

22 <sup>14</sup> German Consortium for Translational Cancer Research DKFZ, Heidelberg, Germany

23 <sup>15</sup> National Center for Tumor Diseases (NCT), Dresden, Germany

24 **Running title:** Synthetic leukemia data with generative AI

25 **Key words:** acute myeloid leukemia, AML, synthetic data, generative model, artificial intelligence

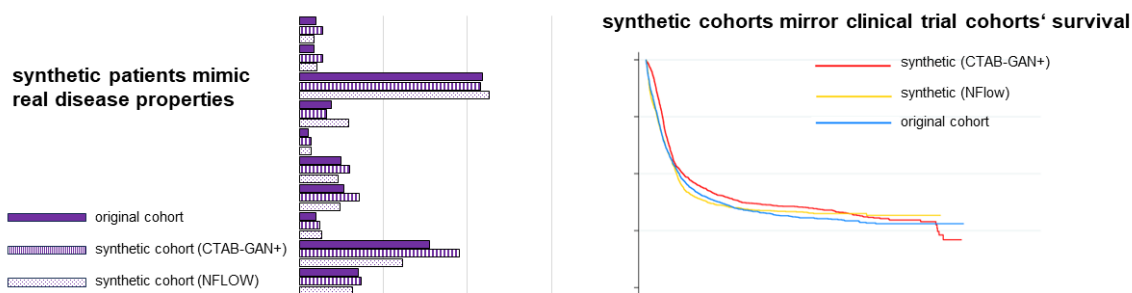
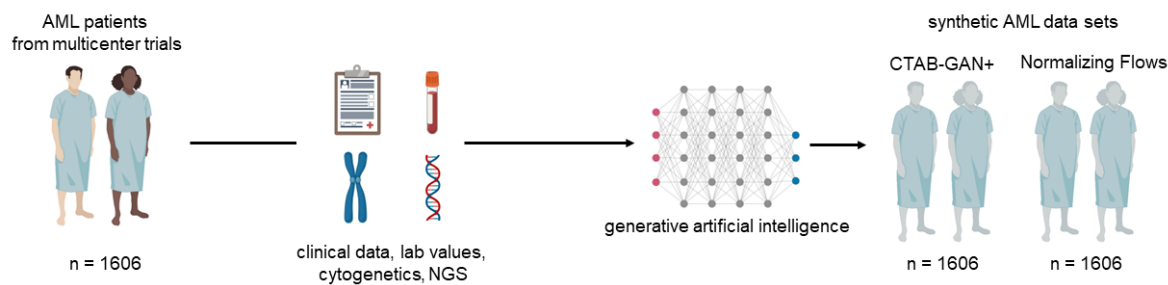
26 **Correspondence:** Jan-Niklas Eckardt, MD, MSc; Department of Internal Medicine I, University  
27 Hospital Carl Gustav Carus, Technical University Dresden and Else-Kröner-Fresenius Center for Digital  
28 Health, Technical University Dresden, Fetscherstraße 74, 01307 Dresden Germany; phone: +49 351 458  
29 11542; e-mail: [jan-niklas.eckardt@uniklinikum-dresden.de](mailto:jan-niklas.eckardt@uniklinikum-dresden.de).

30

31 Word count, abstract: 194; word count, main text: 3876, figures/tables: 6; supplementary figures/tables:  
32 10; references: 39

33 **NOTE:** This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

## 34 Graphical Abstract



35

## 36 Abstract

37 Clinical research relies on high-quality patient data, however, obtaining big data sets is costly and access  
38 to existing data is often hindered by privacy and regulatory concerns. Synthetic data generation holds  
39 the promise of effectively bypassing these boundaries allowing for simplified data accessibility and the  
40 prospect of synthetic control cohorts. We employed two different methodologies of generative artificial  
41 intelligence – CTAB-GAN+ and normalizing flows (NFlow) – to synthesize patient data derived from  
42 1606 patients with acute myeloid leukemia, a heterogeneous hematological malignancy, that were  
43 treated within four multicenter clinical trials. Both generative models accurately captured distributions  
44 of demographic, laboratory, molecular and cytogenetic variables, as well as patient outcomes yielding  
45 high performance scores regarding fidelity and usability of both synthetic cohorts (n=1606 each).  
46 Survival analysis demonstrated close resemblance of survival curves between original and synthetic  
47 cohorts. Inter-variable relationships were preserved in univariable outcome analysis enabling  
48 explorative analysis in our synthetic data. Additionally, training sample privacy is safeguarded  
49 mitigating possible patient re-identification, which we quantified using Hamming distances. We provide  
50 not only a proof-of-concept for synthetic data generation in multimodal clinical data for rare diseases,  
51 but also full public access to synthetic data sets to foster further research.

## 52 **Introduction**

53 In the age of big data, the paucity of publicly available medical data sets is often staggering. Despite  
54 extensive data collection efforts, such as The Cancer Genome Atlas(1), the public availability of  
55 comprehensive entity-specific data sets remains largely unsatisfactory. Data sharing is often hindered  
56 by concerns of patient privacy, regulatory aspects, and proprietary interests.(2) These factors do not only  
57 impede progress in medical research but also establish a gatekeeping mechanism that restricts specific  
58 research inquiries to large institutions with access to extensive datasets. Collecting such data sets is a  
59 costly and time-consuming effort and especially later-phase clinical trials usually take years to complete  
60 and require millions in funding.(3,4) In particular, this is true for rare diseases, such as acute myeloid  
61 leukemia (AML), which is a genetically heterogenous and highly aggressive hematological malignancy  
62 with so far unsatisfactory patient outcomes despite recent advances in therapy.(5) In addition, the  
63 development of targeted therapies for defined subgroups leads to an increased need for control  
64 groups.(6) To gain insights into such burdensome malignant entities with unmet medical needs, a crowd-  
65 sourcing of data to refine risk stratification efforts and test treatment-related hypothesis is essential. If  
66 machine learning methods are to be deployed in such data sets, the size of available diverse training data  
67 is paramount for model robustness. Generative models, especially generative adversarial neural  
68 networks (GANs)(7), have exhibited remarkable capabilities in image generation(8), but can also  
69 effectively generate synthetic non-image data. The unique properties of generative artificial intelligence  
70 (AI) yield the prospect of synthesizing data based on real patients, which can be distributed at will since,  
71 ideally, synthetic data only mimics real patient data alleviating concerns of privacy. In this scenario, the  
72 synthetic data itself should preserve the biological characteristics of the disease under investigation to  
73 make inferences to real-world applications possible. At the same time, synthetic data should safeguard  
74 privacy of the underlying training cohort.

75 In this study, we employ two state-of-the-art technologies of generative modeling on a large training  
76 data set of four pooled multicenter clinical trials including AML patients with comprehensive clinical  
77 and genetic information. We investigate how closely the synthetic data resembles the real trial data  
78 aligning baseline characteristics and patient outcome. Further, we measure privacy conservation in the

79 synthetic data. Additionally, we provide both final fully synthetic data sets comprising 1606 AML  
80 patients each in a publicly accessible repository to foster further research into this devastating disease.

81

## 82 **Methods**

### 83 *Patient data*

84 Multimodal clinical, laboratory, and genetic data (Table S1) were obtained from 1606 patients with non-  
85 M3 AML that were treated within previously conducted multicentric prospective clinical trials of the  
86 German Study Alliance Leukemia (SAL; AML96 [NCT00180115](9), AML2003 [NCT00180102](10),  
87 AML60+ [NCT00180167](11), and SORAML [NCT00893373](12)). Table S2 shows an overview of  
88 trial protocols. Eligibility was determined upon diagnosis of AML, age  $\geq 18$  years, and curative treatment  
89 intent. All patients gave their written informed consent according to the revised Declaration of  
90 Helsinki.(13) All studies were previously approved by the Institutional Review Board of the Technical  
91 University Dresden. Complete remission (CR), event-free survival (EFS), and overall survival (OS)  
92 were defined according to the revised ELN criteria.(14) Biomaterial was obtained from bone marrow  
93 aspirates or peripheral blood prior to treatment initiation. Next-Generation Sequencing (NGS) was  
94 performed using the TruSight Myeloid Sequencing Panel (Illumina, San Diego, CA, USA). Pooled  
95 samples were sequenced paired-end and a 5% variant allele frequency (VAF) mutation calling cut-off  
96 was used with human genome build HG19 as a reference as previously described in detail.(15)  
97 Additionally, high resolution fragment analysis for *FLT3*-ITD(16), *NPM1*(17), and *CEBPA*(18) was  
98 performed as described previously. For cytogenetics, standard techniques for chromosome banding and  
99 fluorescence-in-situ-hybridization (FISH) were used.

100

### 101 *Generative models*

102 In our study, we used two state-of-the-art generative models exhibiting two fundamentally different  
103 concepts of data generation:

104 i) CTAB-GAN+(19) builds upon the Generative Adversarial Network (GAN)(20) architecture,  
105 consisting of two interlinked neural networks - the generator and the discriminator. These are jointly  
106 trained in an adversarial manner. The generator's goal is to produce synthetic data that appears realistic,  
107 starting from random noise. In parallel, the discriminator seeks to differentiate between real and  
108 synthetic samples created by the generator. The training continues until the discriminator is no longer  
109 able to reliably distinguish real data from synthetic, indicating that the generator has successfully  
110 approximated the distribution of the real data.

111 ii) Normalizing Flows (NFlow)(21) presents an alternative approach for synthesizing data from complex  
112 distributions. This comprises a sequence of invertible transformations, starting from a simple base  
113 distribution. Each transformation, or 'flow', gradually modifies this base distribution into a more  
114 complex one that better mirrors the actual data. Importantly, these transformations are stackable,  
115 meaning they can be applied successively to incrementally increase the complexity of the modeled  
116 distribution. All parameters defining these flows are learned directly from the data, allowing the model  
117 to accurately capture the underlying data distribution. Note, that we used a modification of NFlow for  
118 survival data provided by the Synthcity(22) software framework.

119 No imputation of missing data was performed in the original data set, thus both final synthetic data sets  
120 also contain missing data to adequately represent real-world conditions. Hyperparameter tuning was  
121 performed using the Optuna framework allowing both generative models to capture the best possible  
122 representation of the original data. Afterwards, we trained each model with five different random seeds  
123 and sampled from it three times, which generated 15 synthetic datasets for each model. Results are  
124 reported for each highest-performing synthetic data set, respectively.

125

### 126 ***Evaluation of synthetic data performance***

127 To assess the fidelity and usability of synthetic data, previously proposed evaluation metrics were used  
128 to provide a comprehensive overview of model performance. In particular, Basic Statistical Measure,  
129 Regularized Support Coverage, and Log-transformed Correlation Score were used to evaluate the  
130 fidelity of the data in general via our implementation based on the descriptions by Chundawat et al.(23).

131 The second set of metrics – Kaplan-Meier-Divergence, Optimism and Short-Sightedness - was  
132 previously introduced by Norcliffe et al.(24) for synthetic survival data, and implemented in  
133 Synthcity(22). For improved comparability, performance metrics were normalized on a scale from 0  
134 (inadequate representation of original data) to 1 (optimal representation). An overview of the underlying  
135 methodologies of these metrics is provided in Table S3. For detailed information, we refer the interested  
136 reader to the original publications.(23,24)

137

### 138 *Assessment of privacy conservation*

139 To assess potential privacy implications of synthetic data, we customized the method proposed by  
140 Platzer and Reutterer(25) to accommodate for smaller sample sizes. We partitioned the original training  
141 data (80% of total) into four subsets, matching the size of the test dataset (20%) for balanced  
142 comparisons (Fig. S1). Calculations were performed using Hamming distance(26) for categorical  
143 features. Numerical variables were binned (n=10 bins each) and thereby categorized to enable Hamming  
144 distance calculations. Given the nature of the Hamming distance metric, the average minimum distance  
145 effectively denotes the number of variables that would need to be altered for a synthetic patient to match  
146 a real patient. We compared the average distances of the synthetic data to the training (syn → train) and  
147 test sets (syn → test). The relationship between both can be expressed as a coefficient for each synthetic  
148 data set compared to training and test set:

$$149 \quad \text{privacy leakage coefficient} = \frac{\text{syn} \rightarrow \text{test}}{\text{syn} \rightarrow \text{train}} - 1$$

150 By analyzing whether the synthetic data is closer to the training set compared to the test set, we can  
151 assess whether the synthetic data is overly representative of the training data, thereby posing potential  
152 privacy concerns. If the average distances from the synthetic data to the training and test data are equally  
153 small, the privacy leakage coefficient will also be small. The lower the privacy leakage coefficient, the  
154 lower the likelihood of re-identification for patients in the training set. We assumed that values above  
155 0.05 signal potential privacy breaches, as they suggest the synthetic data is substantially closer to the  
156 training set than to the test set. Conversely, values below 0.05 denote a favorable privacy safeguard,

157 signaling similar distances between the training and test sets. Additionally, the number of exact subject  
158 matches between the synthetic and original cohorts was determined.

159

### 160 *Statistical analysis*

161 Pairwise analyses were conducted between the original and both synthetic data sets. Normality was  
162 assessed using the Shapiro-Wilk test. If the assumption of normality was met, continuous variables  
163 between two samples were analyzed using the two-sided unpaired t-test. If the assumption of normality  
164 was violated, continuous variables between two samples were analyzed using the Wilcoxon rank sum  
165 (syn. Mann-Whitney) test. Fisher's exact test was used to compare categorical variables. Univariate  
166 analyses for binary outcomes (CR rate) were carried out via logistic regression to obtain odds ratios  
167 (OR) and 95% confidence intervals (95%-CI). Time-to-event analyses (EFS, OS) were carried out using  
168 Cox proportional hazard models to obtain hazard ratios (HR) and 95%-CI. Kaplan-Meier analyses were  
169 performed for time-to-event data (EFS, OS) and corresponding log-rank tests are reported. Median  
170 follow-up time was calculated using the reverse Kaplan-Meier method.(27) All tests were carried out as  
171 two-sided tests. Statistical significance was determined using a significance level  $\alpha$  of 0.05. Statistical  
172 analysis was performed using STATA BE 18.0 (Stata Corp, College Station, TX, USA).

173

### 174 *Data availability*

175 The final synthetic data sets generated and analyzed for the purpose of this study are publicly available  
176 at <https://zenodo.org/record/8334265>

177

### 178 *Code availability*

179 The code generated for the purpose of this study is publicly available at  
180 [https://github.com/waldemar93/synthetic\\_data\\_pipeline](https://github.com/waldemar93/synthetic_data_pipeline)

181



## 182 **Results**

### 183 **Synthetic cohorts generated by CTAB-GAN+ and NFlow score highly in fidelity metrics**

184 We generated equally sized data sets of  $n=1606$  synthetic patients with each generative model to  
185 compare patient variables to the original cohort. The fidelity of synthetic data was assessed with three  
186 previously proposed performance metrics scaled from 0 (inadequate representation) to 1 (optimal  
187 representation). First, the distribution of each individual variable was compared between original and  
188 synthetic data again yielding high scores for both models (Regularized Support Coverage(23) for  
189 CTAB-GAN+: 0.95 and NFlow: 0.97). Second, continuous numerical variables were assessed by  
190 comparing mean, median, and standard deviation between original and synthetic data per variable (Basic  
191 Statistical Measure(23)) showing high scores for both CTAB-GAN+ (0.91) and NFlow (0.92). Third,  
192 regarding accurate representations of inter-variable correlations, CTAB-GAN+ and NFlow achieved a  
193 Log-Transformed Correlation Score(23) of 0.75 and 0.74, respectively. An overview of performance  
194 metrics is provided in Tab. S4 (usability; survival metrics are reported with survival analysis).

195

### 196 **Synthetic clinical and genetic patient characteristics closely mimic those of real patients**

197 Baseline patient characteristics compared between real and synthetic patients are shown in Table 1. It  
198 has to be noted that given the large sample sizes (three groups with  $n=1606$  each), even small effect  
199 sizes yield statistically significant differences. For instance, median age in the original cohort was 56  
200 years, while synthetic patients generated by CTAB-GAN+ had a slightly younger median age of 53  
201 years ( $p=0.0001$ ), whereas NFlow-generated patients had a slightly older median age of 58 years  
202 ( $p=0.039$ ). Sex distribution did not differ between NFlow and the original cohort, while CTAB-GAN+  
203 generated more males than females (NFlow: 56.2% vs. 43.8%; original: 52.2% vs. 47.8%;  $p=0.023$ ).  
204 The rates of *de novo*, secondary, and therapy-associated AML did not differ significantly for CTAB-  
205 GAN+ generated patients, while NFlow generated fewer *de novo* and more therapy-associated AML  
206 patients compared to the original cohort. Hemoglobin levels and platelet count did not differ  
207 significantly between the original and the synthetic cohorts, while synthetic patients generated by  
208 CTAB-GAN+ showed a significantly higher median white blood cell count than the original cohort.



209 Fifty molecular and cytogenetic alterations were included in generating synthetic patients. Figure 1  
210 displays the distribution of these alterations across the original and synthetic cohorts (absolute numbers  
211 and *p*-values are provided in Tab. S5). These alterations encompass genes that code for epigenetic  
212 regulators (Fig. 1A), the cohesin complex (Fig. 1B), transcription factors (Fig. 1C), *TP53* and  
213 *Nucleophosmin 1* (Fig. 1D), signaling factors (Fig. 1E), components of the spliceosome (Fig. 1F), and  
214 cytogenetic aberrations with established impact on patient outcome (Fig. 1G). Overall, the rates of  
215 alterations in both synthetic cohorts were in a plausible range with a few deviations from the original  
216 cohort of high statistical significance, such as NFlow-generated frequencies of *BCORLI*, *DNMT3A*,  
217 *PHF6*, and *ZRSR2*, as well as CTAB-GAN+-generated frequencies of *CUX1* and *GATA2* while the  
218 remainder of alterations showed only negligible differences. Aside from the frequency per individual  
219 alteration, the co-occurrences of alterations play an important role in disease biology, which should be  
220 also captured in high-quality synthetic data. Fig. 2 shows the relative differences between the original  
221 cohort and CTAB-GAN+ (Fig. 2A) and NFlow (Fig. 2B) regarding co-occurring mutations. We found  
222 high congruencies for co-occurrences compared to the original cohort, while deviations were commonly  
223 found in alterations that had a low frequency in the original cohort.

224

### 225 **Synthetic cohorts match real patients in outcome and survival analysis**

226 Median follow-up for the original cohort was 89.5 months (95%-CI: 85.5-95.4). The synthetic cohorts  
227 had a median follow-up of 91.3 months (CTAB-GAN+, 95%-CI: 84.8-98.0) and 74.3 months (NFlow,  
228 95%-CI: 70.9-77.4). Table 2 shows a detailed comparison of patient outcome between the original and  
229 both synthetic cohorts. For CR rates, we found no significant differences between the original (70.7%)  
230 and both synthetic cohorts (CTAB-GAN+: 73.7%; NFlow: 69.1%). Median EFS in the original cohort  
231 was 7.2 months while both CTAB-GAN+ with 12.8 months and NFlow with 9.0 months deviated with  
232 high significance. This effect can arguably be attributed to both CR rate and OS being included in  
233 hyperparameter tuning, while EFS was exempt from hyperparameter tuning. Kaplan-Meier analysis  
234 nevertheless showed a plausible representation of the survival curves for both synthetic cohorts  
235 regarding EFS (Fig. 3A). Median OS for the original cohort was 17.5 months while the CTAB-GAN+

236 cohort had a median OS of 19.5 months ( $p < 0.0001$ ) and NFlow of 16.2 months ( $p = 0.055$ ). Kaplan-Meier  
237 analysis (Fig. 3B) showed similar behavior of survival curves as for EFS. This was also evident with  
238 regard to usability metrics for synthetic survival data introduced by Norcliffe et al.(24): We found both  
239 CTAB-GAN+ and NFlow to score high in our test set with normalized performance results (+1 is  
240 optimal representation, 0 is inadequate representation, Tab. S4). Kaplan-Meier-Divergence, i.e. the  
241 degree to which survival curves of synthetic and real data differ, was low for both synthetic data sets  
242 (CTAB-GAN+: 0.97, NFlow: 0.98). Neither model showed overt optimism or overt pessimism in  
243 representing survival data (CTAB-GAN+: 0.98, NFlow: 0.99). For both EFS and OS, the curve of  
244 CTAB-GAN+ showed no stabilization of survival rates towards the end of the follow-up period in  
245 comparison to the curve of the original cohort while NFlow tends to censor a higher rate of patients after  
246 passing the 2-year follow-up mark. Nonetheless, Short-sightedness, i.e. failure to predict beyond a  
247 certain time point, was also low for both models, however slightly favoring CTAB-GAN+ over NFlow  
248 (CTAB-GAN+: 0.99, NFlow: 0.93) arguably corresponding to the censoring tendency of NFlow.

249

## 250 **Synthetic data captures risk associations of individual variables for explorative analyses**

251 In order to be useful for explorative analyses, synthetic data needs to recapitulate risk associations of  
252 individual variables. The ELN2022 recommendations represent one of the most widely used guidelines  
253 for risk stratification.(14) Hence, previously established markers of favorable (normal karyotype,  
254  $t(8;21)$ ,  $inv(16)$  or  $t(16;16)$  mutations of *NPM1*, *CEBPA*-bZIP in frame mutations), intermediate risk  
255 (*FLT3*-ITD,  $t(9;11)$ ), or adverse risk (complex karyotype, -5,  $del(5q)$ , -7, -17, mutations of *TP53*,  
256 *RUNX1*, *ASXLI*), and age were evaluated using univariable analyses per cohort for their impact on  
257 achievement of CR, EFS, and OS. All effects for achievement of CR, EFS, and OS showed the same  
258 directionality – favorable affects in the original cohort were also favorable in synthetic cohorts and *vice*  
259 *versa* – and significance – effects that were significant in the original cohort were also significant in  
260 synthetic cohorts and *vice versa* (except for  $del(5q)$  being significantly associated with failure to achieve  
261 CR in the original cohort while this effect turned out to be non-significant in the NFlow-generated  
262 cohort). Importantly, no inverse effects – a variable that would be favorable in the original cohort would

263 be adverse in a synthetic cohort or *vice versa* – were observed. Detailed outcomes per variable are  
264 reported for CR (Tab. S6), EFS (Tab. S7), and OS (Tab. S8).

265

## 266 **Synthetically generated cohorts safeguard real patient data and prohibit re-identification**

267 Privacy conservation was measured by: i) number of exact matches between original and synthetic  
268 cohorts, ii) a privacy leakage coefficient based on Hamming distance, and iii) absolute Hamming  
269 distances showing the number of variables to be altered per synthetic patient to match a real patient.  
270 First, for both synthetic data sets the number of exact matches compared to the original cohort was zero.  
271 Second, the average minimum distances compared between datapoints in training and test sets were  
272 similar for the original cohort, as well as synthetic data from both CTAB-GAN+ and NFlow (Tab. 3).  
273 The privacy leakage coefficient – the quotient of Hamming distances between synthetic to test divided  
274 by synthetic to training data where small values ( $<0.05$ ) indicate a small difference between the distances  
275 of synthetic data to training and test data, and therefore, indicate no privacy breach – was very low for  
276 both CTAB-GAN+ and NFlow (Tab. 3). This signals a low likelihood of re-identification for both  
277 synthetic datasets. Third, the median number of variables that would have to be altered to assign a  
278 synthetic patient to a training set patient was nine for both CTAB-GAN+ and NFlow.

279

## 280 **Discussion**

281 Synthetic data provide an attractive solution to circumvent issues in current standards of data collection  
282 and sharing. These issues encompass first and foremost the time- and cost-intensive data collection  
283 process that usually involves enrollment of patients in prospective clinical trials presenting ever-  
284 increasing costs both regarding funding and time until completion, as well as ethical concerns inherent  
285 in clinical research with human subjects.(3,4) The prospect of using synthetic data as a novel kind of  
286 control group in prospective trials while effectively alleviating the need to enroll a larger number of  
287 patients and cutting costs bears the question of how closely such synthetic control arms match real-  
288 world cohorts. We used two generative AI technologies, a state-of-the-art GAN, CTAB-GAN+, and  
289 NFlow, to mimic the distribution of patient variables from four different previously conducted

290 prospective multicenter trials including a total of 1606 patients with AML. Both models demonstrated  
291 high performance in previously established evaluation metrics that assess fidelity and usability of  
292 synthetic tabular data.(23,24) The comparison of distributions per variable between original and real  
293 data further showed close resemblances. Notably, even for statistically significant deviations from the  
294 original cohort, differences in effect sizes (e.g. age difference, difference in rates of occurrence for  
295 genetic alterations etc.) were often small. Inherent to hypothesis testing with such large sample sizes,  
296 even clinically irrelevant deviations can yield statistically significant differences. Importantly, inter-  
297 variable relationships were conserved in synthetic data: In univariable analyses both effect direction and  
298 statistical significance was well captured by both generative models effectively enabling explorative  
299 investigations in such data sets.

300 Once real data is obtained, privacy concerns often inhibit public access and thus impede data sharing  
301 and third-party hypothesis testing. Frequently used practices range from de-identifying or anonymizing  
302 data to more advanced computational approaches. De-identification or anonymization (e.g. removing  
303 names and birth dates), as well as adding artificial noise to the original data have recently been proven  
304 to be unsafe in terms of guarding privacy as reidentification attacks can successfully unveil patients'  
305 identity.(28–30) Computational advances in both federated(31) and swarm learning(32) where machine  
306 learning models are trained across multiple locations and only either models or weights are shared rather  
307 than the data itself provide a viable alternative. Nevertheless, these technologies are vulnerable to data  
308 reconstructions, e.g. via data leakage from model gradients.(33–35) Inherent to synthetic data generation  
309 in terms of privacy safeguards is a trade-off between usability and privacy where an increase in each  
310 negatively affects the other.(36) Ideally, synthetic data should not be re-identifiable but at the same time  
311 closely match the original distributions. Zero exact matches were observed in our synthetic cohorts.  
312 Additionally, Hamming distances showed that reconstruction of original training samples is highly  
313 unlikely given the number of variables per synthetic patient that would have to be altered in order to  
314 match a training cohort patient.

315 The generation of synthetic data is, as all machine learning models are, fundamentally limited by the  
316 data that the model is trained on. This implies that external users should be aware of the properties of

317 the training data that went into the generation of a synthetic data set in order to either select the right  
318 data set for their research question or *vice versa*, adapt the research question to the available data. It is  
319 therefore important to note, that patients in our trials have all been treated with intensive anthracycline-  
320 based therapy and largely stem from a Middle-European ethnic background. Hence, our generated  
321 synthetic AML data sets may not fully capture features of other populations let alone other treatment  
322 modalities, such as less intensive therapy or targeted agents. The incorporation of these modalities will  
323 be addressed in future works. Since ML models thrive on large and diverse data sets, synthetic data  
324 generation from medical records is caught in a paradoxical loop: Available data is sparse, synthetic data  
325 can potentially accommodate for sparse available real data, synthetic data requires large and diverse sets  
326 of real data to meaningfully represent the population.(37) Therefore, the generation of synthetic data is  
327 likely more robust, if training data from large multicenter cohorts is used. Nonetheless, the availability  
328 of synthetic data promises a democratization of clinical research. In similar efforts, Azizi et al.(38) and  
329 D’Amico et al.(39) explored synthetic data generation in cancer. Azizi et al.(38) used data from a  
330 previously conducted clinical trial in colorectal cancer to generate synthetic data using conditional  
331 decision trees. Focusing on myelodysplastic neoplasms (MDS), D’Amico et al.(39) used a conditional  
332 Wasserstein tabular GAN to generate synthetic MDS patients from the GenoMed4All database. Both  
333 groups conclude the feasibility of either method to generate synthetic data that closely resemble the  
334 original data distributions and provide access to their synthetic data. Such studies may alleviate a  
335 common gatekeeping mechanism of costly data collection efforts that are often restricted to large well-  
336 funded medical centers. Further, this also extends to cross-domain applications involving medical data,  
337 e.g. the training of a ML model by a third party that requires large sets of training data.

338 In summary, we demonstrate the feasibility of two different technologies of generative AI to create  
339 synthetic clinical trial data that both closely mimic disease biology and clinical behavior, as well as  
340 conserve the privacy of patients in the training cohort. Generating such large synthetic data sets based  
341 on multicenter clinical trial training data holds the promise of enabling a new kind of clinical research  
342 improving upon data accessibility, while ameliorating current hindrances in data sharing.

343

## 344 **Acknowledgements**

345 We thank all contributing physicians, laboratories, and nurses associated with the German Study  
346 Alliance Leukemia and especially participating patients for their valuable contributions.

347

## 348 **Authorship Contributions**

349 J.-N.E., W.H., M.W., and J.M.M. designed the study. J.-N.E., C.R, U.P., C. M.-T., H.S., C.D.B., C.S.,  
350 K.S-E., M.H., M.K., A.B., C.T., J.S., M.B., and J.M.M. provided patient samples. S.S. and C.T.  
351 performed molecular analysis. W.H. trained generative models. J.-N.E. performed statistical analysis  
352 and wrote the initial draft. All authors had access to all of the data, analyzed the data, provided critical  
353 scientific insights and revised the draft. All authors agreed to the final version of the manuscript and the  
354 decision to submit it for publication.

355

## 356 **Disclosure of Conflicts of Interest**

357 The authors declare no competing interests.

358

## 359 **References**

- 360 1. The Cancer Genome Atlas Program - National Cancer Institute [Internet]. 2018 [cited 2020 Sep  
361 1]. Available from: [https://www.cancer.gov/about-nci/organization/ccg/research/structural-](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga)  
362 [genomics/tcga](https://www.cancer.gov/about-nci/organization/ccg/research/structural-genomics/tcga)
- 363 2. Taitsman JK, Grimm CM, Agrawal S. Protecting Patient Privacy and Data Security. *New England*  
364 *Journal of Medicine*. 2013 Mar 14;368(11):977–9.
- 365 3. Stewart DJ, Stewart AA, Wheatley-Price P, Batist G, Kantarjian HM, Schiller J, et al. The  
366 importance of greater speed in drug development for advanced malignancies. *Cancer Med*. 2018  
367 Mar 30;7(5):1824–36.
- 368 4. Martin L, Hutchens M, Hawkins C, Radnov A. How much do clinical trials cost? *Nature Reviews*  
369 *Drug Discovery*. 2017 Jun 1;16(6):381–2.
- 370 5. Döhner H, Weisdorf DJ, Bloomfield CD. Acute Myeloid Leukemia. *New England Journal of*  
371 *Medicine*. 2015 Sep 17;373(12):1136–52.

- 372 6. Estey E, Othus M, Gale RP. New study-designs to address the clinical complexity of acute myeloid  
373 leukemia. *Leukemia*. 2019 Mar;33(3):567–9.
- 374 7. Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative  
375 Adversarial Networks [Internet]. arXiv; 2014 [cited 2022 Jul 21]. Available from:  
376 <http://arxiv.org/abs/1406.2661>
- 377 8. Kazemina S, Baur C, Kuijper A, van Ginneken B, Navab N, Albarqouni S, et al. GANs for medical  
378 image analysis. *Artificial Intelligence in Medicine*. 2020 Sep 1;109:101938.
- 379 9. Röllig C, Thiede C, Gramatzki M, Aulitzky W, Bodenstein H, Bornhäuser M, et al. A novel  
380 prognostic model in elderly patients with acute myeloid leukemia: results of 909 patients  
381 entered into the prospective AML96 trial. *Blood*. 2010 Aug 12;116(6):971–8.
- 382 10. Schaich M, Parmentier S, Kramer M, Illmer T, Stölzel F, Röllig C, et al. High-dose cytarabine  
383 consolidation with or without additional amsacrine and mitoxantrone in acute myeloid leukemia:  
384 results of the prospective randomized AML2003 trial. *J Clin Oncol*. 2013 Jun 10;31(17):2094–102.
- 385 11. Röllig C, Kramer M, Gabrecht M, Hänel M, Herbst R, Kaiser U, et al. Intermediate-dose cytarabine  
386 plus mitoxantrone versus standard-dose cytarabine plus daunorubicin for acute myeloid  
387 leukemia in elderly patients. *Ann Oncol*. 2018 01;29(4):973–8.
- 388 12. Röllig C, Serve H, Hüttmann A, Noppeney R, Müller-Tidow C, Krug U, et al. Addition of sorafenib  
389 versus placebo to standard therapy in patients aged 60 years or younger with newly diagnosed  
390 acute myeloid leukaemia (SORAML): a multicentre, phase 2, randomised controlled trial. *Lancet*  
391 *Oncol*. 2015 Dec;16(16):1691–9.
- 392 13. World Medical Association. World Medical Association Declaration of Helsinki: Ethical Principles  
393 for Medical Research Involving Human Subjects. *JAMA*. 2013 Nov 27;310(20):2191–4.
- 394 14. Döhner H, Wei AH, Appelbaum FR, Craddock C, DiNardo CD, Dombret H, et al. Diagnosis and  
395 Management of AML in Adults: 2022 ELN Recommendations from an International Expert Panel.  
396 *Blood*. 2022 Jul 7;blood.2022016867.
- 397 15. Stasik S, Schuster C, Ortlepp C, Platzbecker U, Bornhäuser M, Schetelig J, et al. An optimized  
398 targeted Next-Generation Sequencing approach for sensitive detection of single nucleotide  
399 variants. *Biomol Detect Quantif*. 2018 May;15:6–12.
- 400 16. Thiede C, Steudel C, Mohr B, Schaich M, Schäkel U, Platzbecker U, et al. Analysis of FLT3-  
401 activating mutations in 979 patients with acute myelogenous leukemia: association with FAB  
402 subtypes and identification of subgroups with poor prognosis. *Blood*. 2002 Jun 15;99(12):4326–  
403 35.
- 404 17. Thiede C, Koch S, Creutzig E, Steudel C, Illmer T, Schaich M, et al. Prevalence and prognostic  
405 impact of NPM1 mutations in 1485 adult patients with acute myeloid leukemia (AML). *Blood*.  
406 2006 May 15;107(10):4011–20.
- 407 18. Taube F, Georgi JA, Kramer M, Stasik S, Middeke JM, Röllig C, et al. CEBPA Mutations in 4708  
408 Patients with Acute Myeloid Leukemia - Differential Impact of bZIP and TAD Mutations on  
409 Outcome. *Blood*. 2021 Jul 28;blood.2020009680.
- 410 19. Zhao Z, Kunar A, Birke R, Chen LY. CTAB-GAN+: Enhancing Tabular Data Synthesis [Internet].  
411 arXiv; 2022 [cited 2023 Jul 24]. Available from: <http://arxiv.org/abs/2204.00401>



- 412 20. Goodfellow IJ, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, et al. Generative  
413 Adversarial Networks. arXiv:14062661 [cs, stat] [Internet]. 2014 Jun 10 [cited 2021 May 27];  
414 Available from: <http://arxiv.org/abs/1406.2661>
- 415 21. Papamakarios G, Nalisnick E, Rezende DJ, Mohamed S, Lakshminarayanan B. Normalizing Flows  
416 for Probabilistic Modeling and Inference [Internet]. arXiv; 2021 [cited 2023 Jul 24]. Available  
417 from: <http://arxiv.org/abs/1912.02762>
- 418 22. Qian Z, Cebere BC, van der Schaar M. Synthcity: facilitating innovative use cases of synthetic data  
419 in different data modalities [Internet]. arXiv; 2023 [cited 2023 Jul 24]. Available from:  
420 <http://arxiv.org/abs/2301.07573>
- 421 23. Chundawat VS, Tarun AK, Mandal M, Lahoti M, Narang P. TabSynDex: A Universal Metric for  
422 Robust Evaluation of Synthetic Tabular Data [Internet]. arXiv; 2022 [cited 2023 Jul 24]. Available  
423 from: <http://arxiv.org/abs/2207.05295>
- 424 24. Norcliffe A, Cebere B, Imrie F, Lio P, van der Schaar M. SurvivalGAN: Generating Time-to-Event  
425 Data for Survival Analysis [Internet]. arXiv; 2023 [cited 2023 Aug 3]. Available from:  
426 <http://arxiv.org/abs/2302.12749>
- 427 25. Platzer M, Reutterer T. Holdout-Based Fidelity and Privacy Assessment of Mixed-Type Synthetic  
428 Data [Internet]. arXiv; 2021 [cited 2023 Aug 10]. Available from: <http://arxiv.org/abs/2104.00635>
- 429 26. Hamming RW. Error detecting and error correcting codes. *The Bell System Technical Journal*.  
430 1950 Apr;29(2):147–60.
- 431 27. Shuster JJ. Median follow-up in clinical trials. *J Clin Oncol*. 1991 Jan;9(1):191–2.
- 432 28. Emam KE, Jonker E, Arbuckle L, Malin B. A Systematic Review of Re-Identification Attacks on  
433 Health Data. *PLOS ONE*. 2011 Dec 2;6(12):e28071.
- 434 29. Ursin G, Sen S, Mottu JM, Nygård M. Protecting Privacy in Large Datasets-First We Assess the  
435 Risk; Then We Fuzzy the Data. *Cancer Epidemiol Biomarkers Prev*. 2017 Aug 1;26(8):1219–24.
- 436 30. Sweeney L, Yoo JS, Perovich L, Boronow KE, Brown P, Brody JG. Re-identification Risks in HIPAA  
437 Safe Harbor Data: A study of data from one environmental health study. *Technol Sci*.  
438 2017;2017:2017082801.
- 439 31. Rieke N, Hancox J, Li W, Milletari F, Roth HR, Albarqouni S, et al. The future of digital health with  
440 federated learning. *npj Digit Med*. 2020 Sep 14;3(1):1–7.
- 441 32. Warnat-Herresthal S, Schultze H, Shastry KL, Manamohan S, Mukherjee S, Garg V, et al. Swarm  
442 Learning for decentralized and confidential clinical machine learning. *Nature*. 2021  
443 Jun;594(7862):265–70.
- 444 33. Melis L, Song C, De Cristofaro E, Shmatikov V. Exploiting Unintended Feature Leakage in  
445 Collaborative Learning [Internet]. arXiv; 2018 [cited 2023 Jul 10]. Available from:  
446 <http://arxiv.org/abs/1805.04049>
- 447 34. Zhu L, Liu Z, Han S. Deep Leakage from Gradients [Internet]. arXiv; 2019 [cited 2023 Jul 10].  
448 Available from: <http://arxiv.org/abs/1906.08935>

- 449 35. Boenisch F, Dziedzic A, Schuster R, Shamsabadi AS, Shumailov I, Papernot N. When the Curious  
450 Abandon Honesty: Federated Learning Is Not Private [Internet]. arXiv; 2023 [cited 2023 Jul 10].  
451 Available from: <http://arxiv.org/abs/2112.02918>
- 452 36. Rajotte JF, Bergen R, Buckeridge DL, El Emam K, Ng R, Strome E. Synthetic data as an enabler for  
453 machine learning applications in medicine. *iScience*. 2022 Nov 18;25(11):105331.
- 454 37. Chen RJ, Lu MY, Chen TY, Williamson DFK, Mahmood F. Synthetic data in machine learning for  
455 medicine and healthcare. *Nat Biomed Eng*. 2021 Jun;5(6):493–7.
- 456 38. Azizi Z, Zheng C, Mosquera L, Pilote L, Emam KE. Can synthetic data be a proxy for real clinical  
457 trial data? A validation study. *BMJ Open*. 2021 Apr 1;11(4):e043497.
- 458 39. D’Amico S, Dall’Olio D, Sala C, Dall’Olio L, Sauta E, Zampini M, et al. Synthetic Data Generation by  
459 Artificial Intelligence to Accelerate Research and Precision Medicine in Hematology. *JCO Clinical  
460 Cancer Informatics*. 2023 Jul;(7):e2300021.
- 461

## 462 Tables

clinical data	original cohort	CTAB-GAN+	<i>p</i>	NFlow	<i>p</i>
number of patients	1606	1606		1606	
age, median (IQR)	56 (44 - 65)	53 (42 - 64)	<b>0.0001</b>	58 (47 – 66)	<b>0.039</b>
sex, n (%)			<b>0.023</b>		0.672
female	768 (47.8)	703 (43.8)		781 (48.6)	
male	838 (52.2)	903 (56.2)		825 (51.4)	
AML status, n (%)					
de novo	1339 (83.4)	1339 (83.4)	1.000	1250 (77.8)	<b>0.041</b>
secondary	195 (12.1)	193 (12.0)	0.914	200 (12.5)	0.554
therapy-associated	54 (3.4)	57 (3.5)	0.847	83 (5.2)	<b>0.007</b>
extramedullary disease, n (%)	224 (13.9)	228 (14.2)	0.409	279 (17.4)	<b>0.003</b>
<b>laboratory values</b>					
WBC, median (IQR) in GPt/l	19.5 (4.5 - 53.4)	27.0 (8.3 - 69.6)	<b>&lt;0.0001</b>	14.4 (5.8 – 55.3)	0.832
Hb, median (IQR) in mmol/l	5.9 (5.0 - 8.6)	5.8 (5.0 - 7.0)	0.949	5.9 (5.2 – 6.8)	0.988
Plt, median (IQR) in GPt/l	50.0 (27.0 – 94.0)	49.7 (31.0 - 93.4)	0.073	48.0 (26.2 – 94.5)	0.405

463 **Table 1 Distribution of baseline characteristics between the original and synthetic cohort.** Boldface

464 indicates statistical significance ( $p < 0.05$ ).  $p$ -values are calculated using two-sample comparisons

465 between each of the synthetic cohorts and the baseline cohort for reference. Abbreviations: Hb:

466 hemoglobin; IQR: interquartile range; n: number; Plt: platelet count; WBC: white blood cell count.

467

	original cohort	CTAB-GAN+	NFlow
CR after induction therapy, n (%)	1135 (70.7)	1184 (73.7)	1110 (69.1)
OR	2.41	2.81	2.24
[95%-CI]	[2.16 – 2.68]	[2.51 – 3.14]	[2.01 – 2.49]
$p$ -value		0.059	0.356
median EFS, months (IQR)	7.2 (6.5 – 8.0)	12.8 (11.8 – 14.1)	9.0 (8.3 – 9.7)
HR	1.36	0.74	0.87
[95%-CI]	[1.25 – 1.47]	[0.68 – 0.80]	[0.80 – 0.94]
$p$ -value		<b>&lt;0.0001</b>	<b>&lt;0.0001</b>
median OS, months (IQR)	17.5 (15.7 – 19.2)	19.5 (15.7 – 19.2)	16.2 (15.7 – 19.2)
HR	1.14	0.88	1.00
[95%-CI]	[1.04 – 1.24]	[0.81- 0.96]	[0.92 – 1.09]
$p$ -value		<b>&lt;0.0001</b>	0.055

468 **Table 2 Comparison of patient outcomes between the original and synthetic cohort.** Logistic

469 regression and Cox proportional hazard models were used to obtain odds ratios (OR) for achievement

470 of complete remission (CR) and hazard ratio (HR) with corresponding 95%-confidence intervals (95%-

471 CI). Boldface indicates statistical significance ( $p < 0.05$ ).  $p$ -values are calculated using two-sample

472 comparisons between each of the synthetic cohorts and the original cohort for reference. Other

473 abbreviations: n: number.

474

475

476

477

478

	CTAB-GAN+	NFlow	original cohort
<b>absolute Hamming distances</b>			
average min. distance train	8.7034	9.3474	8.2524
average min. distance test	8.8587	9.4117	8.2224
median distance train	9	9	8
median distance test	9	9	8
<b>relative Hamming distances</b>			
privacy leakage coefficient	0.0178	0.0069	

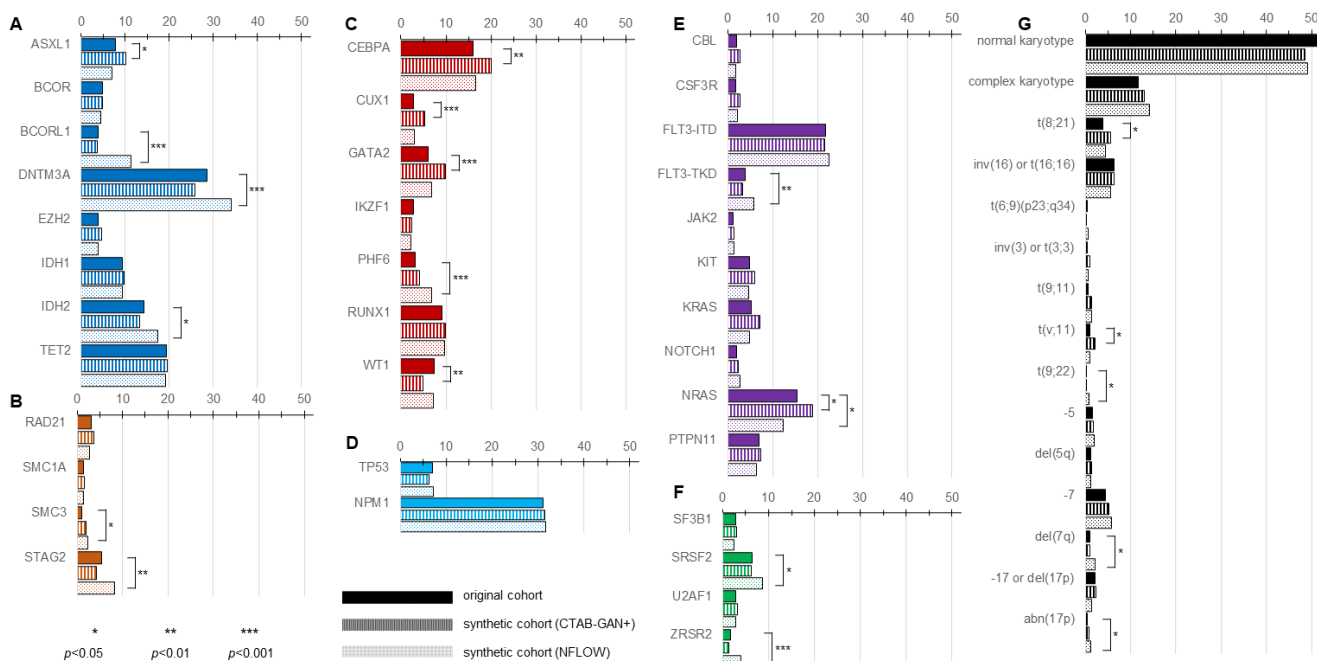
479 **Table 3 Hamming distances for privacy conservation.** Hamming distances were used to measure the  
480 distance between two points within and between equally sized subsets of training (four sets of 20%) and  
481 test data (20%). The median distance represents the number of variables that have to be altered (and  
482 matched exactly) to fit a real patient. A threshold for the privacy leakage coefficient of 0.05 for relative  
483 distances was set where values above 0.05 signal potential privacy breaches. Both synthetic data sets  
484 fell well below the 0.05 threshold signaling larger distances between synthetic and training data, which  
485 make a re-identification of training set patients unlikely.

486

487

## 488 Figures and Figure Legends

Figure 1



489

490

491 **Figure 1 Distribution of molecular and cytogenetic alterations between real and synthetic patients.**

492 50 molecular genetic and cytogenetic alterations were included in generative modeling. Molecular

493 genetics were originally assessed by next-generation sequencing using a targeted myeloid panel

494 including genes that encode for epigenetic regulators (A, dark blue), the cohesion complex (B, orange),

495 transcription factors (C, red), NPM1 and TP53 (D, light blue), signaling factors (E, purple), and the

496 spliceosome (F, green). Cytogenetic aberrations (G, black) were selected based on previously

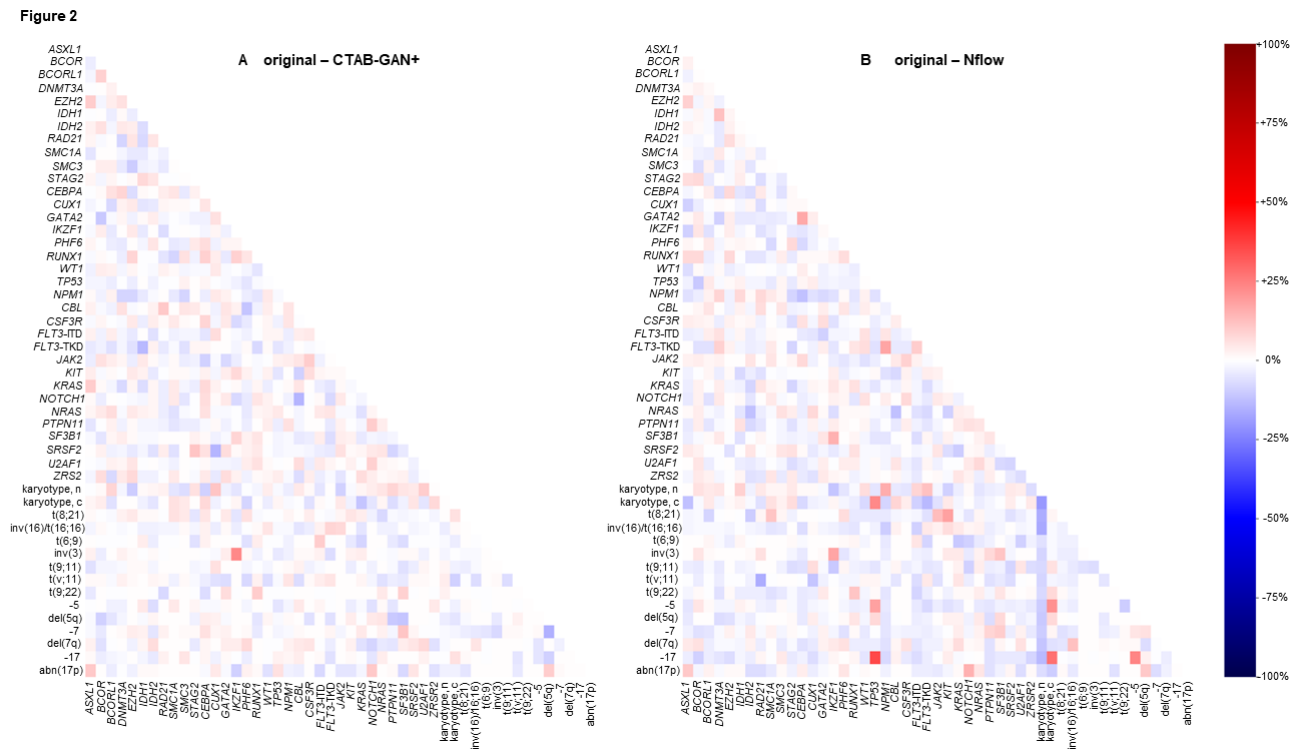
497 demonstrated impact on patient outcomes. Distributions for all variables are denoted as percentages of

498 each respective cohort. Overall, both synthetic cohorts well represented the distribution of alterations in

499 the original cohort with only slight deviations denoted by highly statistically significant ( $p < 0.001$ )

500 differences in *BCORL1*, *DNMT3A*, *PHF6*, and *ZRSR2* for NFlow, as well as *CUX1* and *GATA2* for

501 CTAB-GAN+.



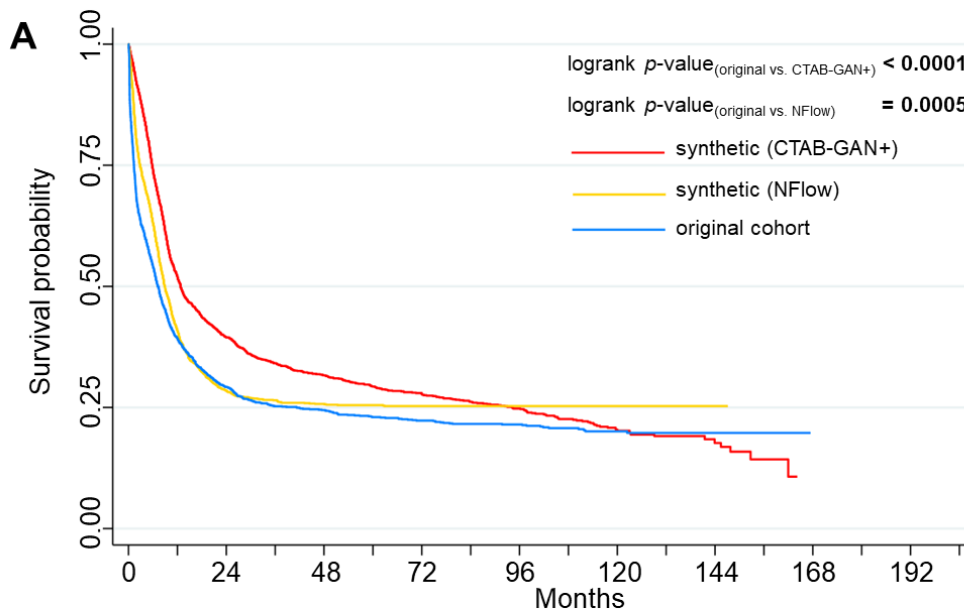
502

503 **Figure 2 Heatmaps for relative differences of genetic associations.** The difference in co-occurrences  
 504 of genetic alterations are plotted. Relative increases (red) or decreases (blue) are displayed on a scale  
 505 from -100% to + 100%. The overlap between the original cohort and CTAB-GAN+ (A), as well as  
 506 original and NFlow (B) showed high congruency. Increases or decreases in co-occurring genetic  
 507 alterations were commonly found to affect alterations with low frequency in the original cohort.

508

**Figure 3**

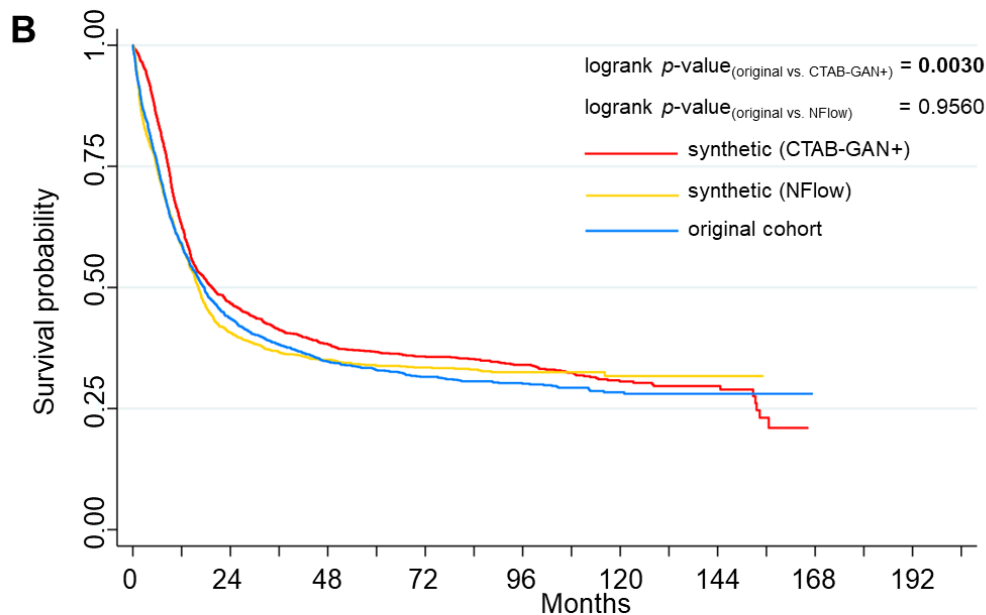
**Event-free survival**



Numbers at risk

real	1606	454	292	238	161	68	27	0
CTAB+	1606	590	407	313	210	97	23	0
NFlow	1606	418	315	227	84	20	1	0

**Overall survival**



Numbers at risk

real	1606	671	398	324	217	96	30	0
CTAB+	1606	699	467	364	257	128	41	0
NFlow	1606	590	407	285	116	34	4	0



510 **Figure 3 Comparison of survival curves between original and synthetic cohorts.** Event-free survival  
511 (EFS) deviated significantly from the original cohort for both synthetic cohorts (A). For the NFlow-  
512 generated cohort, there was no significant deviation from the original distribution for overall survival  
513 (OS), while the CTAB-GAN+-generated cohort again differed significantly (B). Interestingly, while the  
514 survival curve for CTAB-GAN+ displays a plausible curve up until ten years of follow-up, the curve  
515 shows no stabilization of survival rates in the end as the original cohort does. Contrastingly, the survival  
516 curve for NFlow shows an overall plausible course, however, NFlow tends to overtly censor patients  
517 after two years of follow-up.