

1 **A confounder debiasing method for RCT-like comparability enables Machine Learning-**
2 **based personalization of survival benefit in living donor liver transplantation**

3
4 Anirudh Gangadhar*^[a,b], Bima J. Hasjim*^[c], Xun Zhao^[d], Yingji Sun^[b], Joseph Chon^[a], Aman
5 Sidhu^[a], Elmar Jaeckel^[a,f], Nazia Selzner^[a,f], Mark S. Cattral^[a,e], Blayne A. Sayed^[a], Michael
6 Brudno^[g-i], Chris McIntosh^{‡[b,g,h,j-l]}, Mamatha Bhat^{‡[a,b,f,h]}

7
8 *Co-first authors
9 ‡Co-senior authors

10
11 ^[a] Transplant AI initiative, Ajmera Transplant Centre, University Health Network, University of
12 Toronto, ON, Canada

13 ^[b] Toronto General Hospital Research Institute, University Health Network, Toronto, ON,
14 Canada

15 ^[c] Department of Surgery, University of California – Irvine, Orange, California, USA

16 ^[d] McGill University Health Center, Montreal, Quebec, Canada

17 ^[e] Department of Surgery, University of Toronto, Toronto, ON, Canada

18 ^[f] Division of Gastroenterology & Hepatology, Department of Medicine, University of Toronto,
19 Toronto, ON, Canada

20 ^[g] Department of Computer Science, University of Toronto, Toronto, ON, Canada

21 ^[h] Vector Institute, Toronto, ON, Canada

22 ^[i] Princess Margaret Cancer Center, University Health Network, Toronto, ON, Canada

23 ^[j] Joint Department of Medical Imaging, University Health Network, Toronto, ON, Canada

24 ^[k] Department of Medical Biophysics, University of Toronto, Toronto, ON, Canada

25 ^[l] Peter Munk Cardiac Center and Ted Rogers Centre for Heart Research, University Health
26 Network, Toronto, ON, Canada

27
28
29
30 **Corresponding author:**

31 Mamatha Bhat, MD PhD

32 Department of Medicine, Division of Gastroenterology & Hepatology

33 University of Toronto

34 Partnerships & Engagement Lead, Temerty Centre for AI in Research & Education in Medicine
35 (T-CAIREM)

36 Faculty Affiliate, Vector Institute for Artificial Intelligence

37 Email: Mamatha.Bhat@uhn.ca

38

39

40

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71

Abstract

Many clinical questions in medicine cannot be answered through randomized controlled trials (RCTs) due to ethical or feasibility constraints. In such cases, observational data is often the only available resource for evaluating treatment effects. To address this challenge, we have developed Decision Path Similarity Matching (DPSM), a novel machine learning (ML)-based algorithm that simulates RCT-like conditions to debias observational data. In this study, we apply DPSM to the clinical question of living donor liver transplantation (LDLT) versus deceased donor liver transplantation (DDLT), helping to identify which patients benefit most from LDLT. DPSM leverages decision paths from a Random Forest classifier to perform accurate, one-to-one matching between LDLT and DDLT recipients, minimizing confounding while retaining interpretability. Using data from the Scientific Registry of Transplant Recipients (SRTR), including 4,473 LDLT and 68,108 DDLT patients transplanted between 2002 and 2023, we trained independent Random Survival Forest (RSF) models on the matched cohorts to predict post-transplant survival. DPSM successfully reduced confounding associations between the two groups as shown by a decrease in area under the receiver operating characteristic (AUROC) from 0.82 to 0.51. Subsequently, RSF ($C\text{-index}_{ldlt}=0.67$, $C\text{-index}_{ddlt}=0.74$) outperformed the traditional Cox model ($C\text{-index}_{ldlt}=0.57$, $C\text{-index}_{ddlt}=0.65$). The predicted 10-year mean survival gain was 10.3% (SD = 5.7%). In conclusion, DPSM provides an effective approach for creating RCT-like comparability from observational data, enabling personalized survival predictions. By leveraging real-world data where RCTs are impractical, this method offers clinicians a tool for transitioning from population-level evidence to more nuanced, personalization.

72

73 **1. Main**

74 Many clinical questions in medicine cannot be addressed through randomized controlled trials¹
75 (RCTs) due to ethical, logistical, or practical challenges. In liver transplantation (LT), for
76 example, it is not feasible to randomize patients between living donor liver transplantation
77 (LDLT) and deceased donor liver transplantation (DDLT) because of the ethical implications of
78 assigning healthy donors and the urgency for life-saving transplants²⁻⁴. As a result, we often rely
79 on observational data to assess the relative benefits of LDLT and DDLT.

80 Clinically, LDLT offers significant advantages over DDLT, such as reduced waitlist times⁵⁻⁸,
81 improved graft quality^{5,9}, and lower rejection rates¹⁰. Despite these benefits, LDLT remains
82 underutilized, representing only 5% of all LT cases in the United States⁶. Previous studies^{5,10,11}
83 have shown general survival benefits of LDLT compared to DDLT, but they lack sufficient
84 adjustment for confounding factors, making it difficult for clinicians to determine which individual
85 patients would benefit most from LDLT.

86 RCTs are considered the gold standard for assessing intervention effects, but as mentioned,
87 they are not feasible in the LDLT versus DDLT context. This has led to the development of
88 advanced statistical and machine learning methods that can simulate RCT-like conditions using
89 observational data. Propensity Score Matching (PSM)^{12,13} is one such method, but it has
90 limitations. By reducing complex, multi-dimensional covariate space into a single probability
91 score, PSM can fail to balance key variables and interactions, leading to residual confounding
92 and imprecise graft-type effect estimates¹⁴⁻¹⁶.

93 In response to these limitations, we introduce Decision Path Similarity Matching (DPSM), a
94 novel machine learning-based algorithm designed to improve matching by leveraging the
95 decision paths from Random Forest models. Unlike PSM, our method matches patients based
96 on entire decision paths rather than a single probability score. This richer representation
97 captures complex, non-linear relationships between covariates, enabling more precise matching
98 and minimizing confounding. DPSM also allows for explainability by providing per-matched-pair
99 visualization of the key variables driving the decision-making process.

100 After applying DPSM to match LDLT and DDLT patients, we utilize a time-to-event machine
101 learning framework, specifically Random Survival Forest (RSF) models, to predict long-term

102 survival outcomes. To the best of our knowledge, no existing method offers individualized
 103 survival benefit predictions for LDLT versus DDLT based on patient-specific variables, making
 104 this work a significant advancement in the field.

105

106 2. Results

107 2.1. Patient characteristics

108 A total of 72,581 LT recipients were included in the study. DDLTs constituted 93.8% ($n_{adlt} =$
 109 68,108), while LDLTs comprised a much lesser percentage at 6.2%, ($n_{ldlt} = 4,473$).

110 Demographic and clinical study variables for both groups are reported in Table 1. DDLT patients
 111 had higher rates of post-transplant mortality (29.7%) as compared to LDLT (22.3%).

Characteristics	LDLT population (n = 4,473)	DDLT population (n = 68,108)	p-value
Age (years)	51.5 (12.8)	52.9 (10.4)	<0.001
Sex			
Female (%)	48.3	35.6	
Male (%)	51.7	64.4	<0.001
Height (cm)	169.9 (10.3)	172.3 (10.3)	<0.001
Weight (kg)	78.8 (17.7)	87.2 (20.3)	<0.001
BMI	27.2 (5.2)	29.3 (5.9)	<0.001
Blood type (%)			
A	43.3	36.7	<0.001
AB	1.9	5.4	<0.001
B	9.8	13.5	<0.001
O	45.0	44.4	0.375
MELD	14.0 (4.5)	22.4 (8.7)	0.000
Primary etiology (%)			
Alcoholic Cirrhosis	15.1	29.4	<0.001
Autoimmune Hepatitis	4.6	3.5	<0.001
Cryptogenic Cirrhosis	6.9	6.6	0.378
HBV Cirrhosis	1.0	1.7	<0.001
HCV Cirrhosis	15.8	22.1	<0.001
Metabolic Liver Disease	2.0	2.2	0.374
MASH	16.9	15.3	0.007
PSC	19.0	5.2	0.956
PBC	8.9	3.7	<0.001
Other	8.0	7.9	<0.001

112 **Table 1. Clinicodemographic characteristics of LDLT, DDLT recipients.**

113

114 2.2. LDLT-DDLT matching

115 In our study, where the goal is to estimate the survival benefit associated with receiving one
116 type of intervention (LDLT) over another (DDLT), it becomes necessary to minimize
117 confounding associations to ensure that our findings are not subject to bias. To this end, we
118 developed the DPSM matching algorithm (details in Sec. 4.4) that performs optimal one-to-one
119 matching between LDLT and DDLT patients using all study variables listed in Sec. 4.2 (Fig. 1a).
120 Unlike propensity score matching (PSM), which matches patients based on an output probability
121 scalar, DPSM leverages and scores entire decision paths produced by a Random Forest model.
122 This approach not only enhances the accuracy of the matching process but also increases the
123 explainability of the method, providing a more transparent and interpretable framework for
124 understanding factors influencing the matching.

125
126 First, we evaluated the effectiveness of our matching technique. Fig. 1b compares 1-D, 2-D pre-
127 and post-match distributions for 2 key variables: age and MELD score across LDLT, DDLT
128 patients. We observe a high degree of overlap for the matched populations, confirming the
129 success of our method. Originally, DDLT patients had a relatively higher MELD score ($22.4 \pm$
130 8.7) than those that received an LDLT (14.0 ± 4.5). High MELD (>33) patients, generally much
131 sicker, are unable to be matched as they do not possess an LDLT counterpart. All other
132 variables used in the study were also found to exhibit good matching (Fig. S2). Additionally, we
133 sought to understand how matching impacts survival times of LDLT and DDLT patients. Kaplan-
134 Meier analysis estimated LDLT and DDLT median survival times in the original dataset to be
135 219 and 184 months respectively (Fig. 1c). After matching using our DPSM technique, we
136 observed a slight decrease in median survival times for the two groups ($t_{\text{surv, ldlt}} = 215$ months,
137 $t_{\text{surv, ddlt}} = 176$ months) (Fig. 1d). For completeness, survival times after PSM matching were also
138 computed ($t_{\text{surv, ldlt}} = 219$ months, $t_{\text{surv, ddlt}} = 192$ months) (Fig. 1e).

139
140 Next, we quantitatively evaluated the efficacy of our DPSM method by computing the AUROC of
141 a Random Forest classifier trained on the pre- and post-matched datasets (Fig. 1f). The mean
142 AUROC performance dropped significantly from 0.83 on the original dataset to 0.51 after
143 matching, indicating that the model found it challenging to accurately classify patients as LDLT
144 or DDLT, thus demonstrating the success of our matching process in reducing systematic
145 differences between the two groups. Notably, post-match AUROC for DPSM (0.509 ± 0.018)
146 was lower than that achieved using traditional PSM (0.589 ± 0.007), suggesting that our method
147 performed better in achieving balanced and comparable groups. This enables establishment of
148 a clearer causal relationship between graft type and survival outcomes.

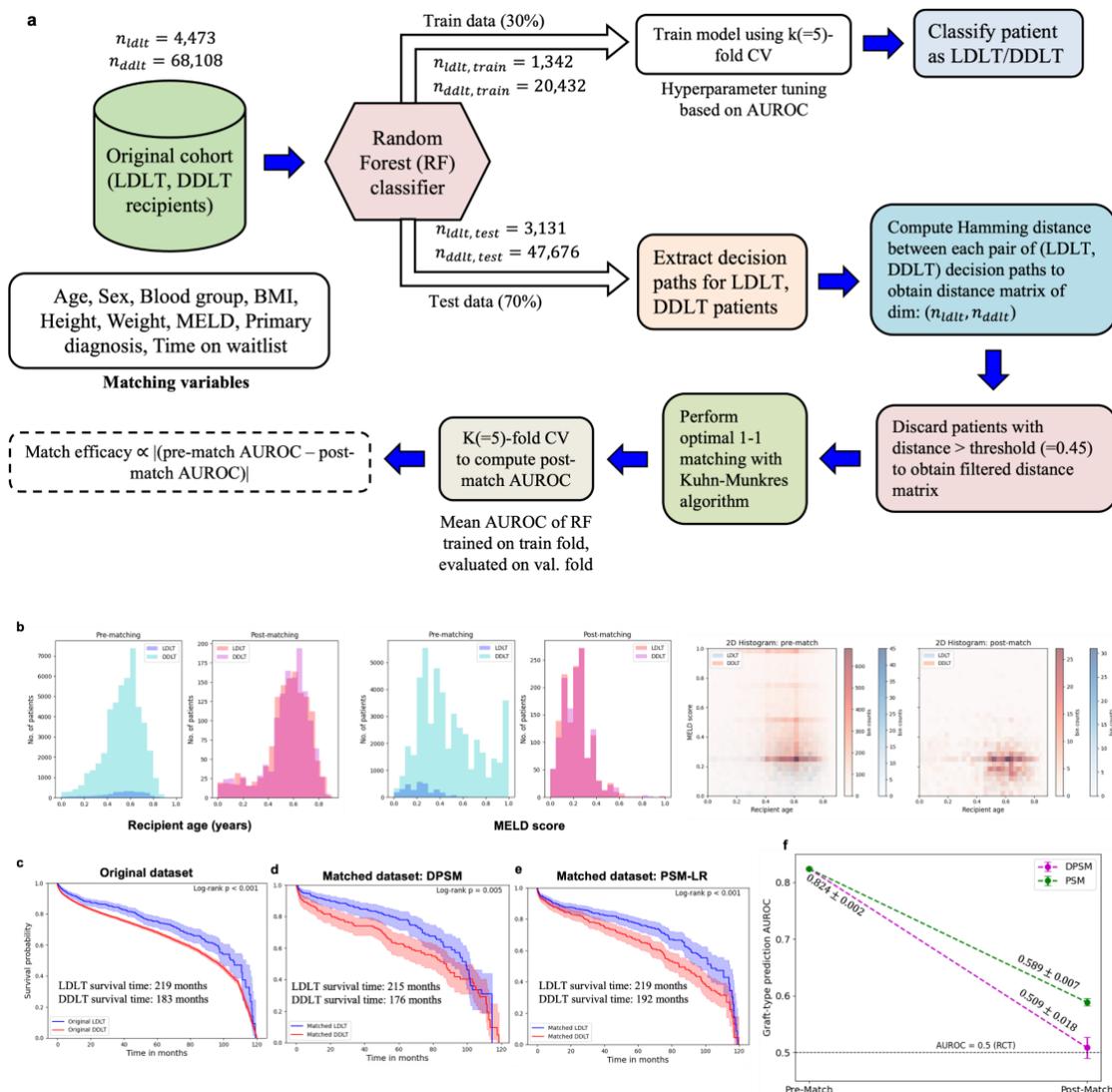


Fig. 1. LDLT-DDLT Matching.

(a) Workflow of our DPSM algorithm; (b) 1-D/2-D distributions (histogram) for recipient age and MELD score variables pre- and post-matching; (c) observed LDLT, DDLT survival in original dataset (pre-match); (d) observed LDLT, DDLT survival post-DPSM-matching; (e) observed LDLT, DDLT survival post-PSM-matching; (f) post-match AUROCs achieved by our DPSM and traditional PSM methods for the graft-type prediction task. The ideal case is AUROC=0.5, which corresponds to an actual RCT. Error bars indicate standard deviations across 10 random samplings.

149

150 A key advantage of our DPSM method is its ability to provide explainable matching, enhancing
 151 the transparency of the model's decision-making process. To illustrate this, we show the 10
 152 most frequently occurring patient variables in the decision paths across all trees of the Random
 153 Forest. This type of interpretability is essential because it transforms a complex, "black-box"
 154 setup into something that can be understood by researchers and clinicians. Fig. 2 shows these
 155 explanations for three randomly selected LDLT-matched DDLT patient pairs, where we observe

156 that MELD score, waitlist time, weight and BMI are the top variables that influence the model
157 deciding whether a patient received an LDLT or DDLT. Other patient characteristics such as
158 age, sex and height were also found to be important predictors. Among etiologies, Alcoholic
159 Cirrhosis, HCV Cirrhosis and PSC were deemed important by the Random Forest model to
160 make graft-type predictions.

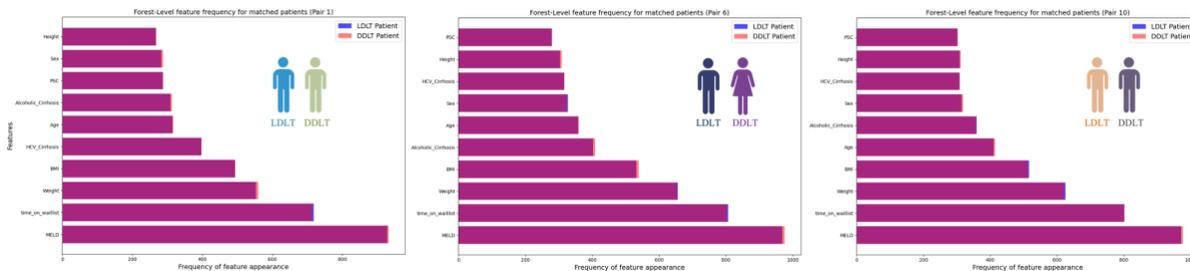


Fig. 2. Frequently occurring patient characteristics across forest-level decision paths.

For a given LDLT-matched DDLT patient pair, we show the 10 most frequently occurring variables across decision paths across all Trees of the Random Forest. A high degree of overlap indicates similarity between the two patients (differing by graft-type) in terms of how the forest makes decisions. This is illustrated for 3 randomly chosen patient pairs.

161

162 2.3. Evaluation of survival model

163 Building on the matched LDLT and DDLT cohorts generated by DPSM, we trained and
164 evaluated survival models to predict patient survival outcome. Methodological details are
165 provided in Sec. 4.5. For this task, we compared the performance of two popular time-to-event
166 models, Random Survival Forest (RSF) and Cox Proportional Hazards (CPH). Two independent
167 models were trained on the LDLT and matched DDLT populations.

168

169 In general, RSF ($C - index_{ldlt} = 0.673$, $C - index_{adlt} = 0.740$) performed better than CPH ($C -$
170 $index_{ldlt} = 0.572$, $C - index_{adlt} = 0.652$) on the C-index and was therefore selected as the
171 model of choice (Fig 3b). The improved performance may be attributed to the former's ability to
172 model non-linear data patterns without making explicit assumptions about underlying
173 distributions. The average Brier score did not exceed 0.14. Additionally, we performed SHAP
174 analysis to understand feature contributions to our outcome of interest, i.e., post-transplant
175 mortality (Fig. 3c, d). Recipient age emerged as the strongest predictor of mortality with older
176 patients at much greater risk. Other important factors were weight, BMI, MELD, height and
177 blood type A. In terms of the primary indication for transplant, PSC, MASH, Alcoholic Cirrhosis
178 and HCV Cirrhosis were all found to be important risk factors. In light of these findings, it
179 becomes important to point out that variables such as age, MELD score, weight and BMI have

180 been identified as key confounders. This is because these variables impact both, graft-type
 181 assignment (Fig. 2) as well as survival outcome.

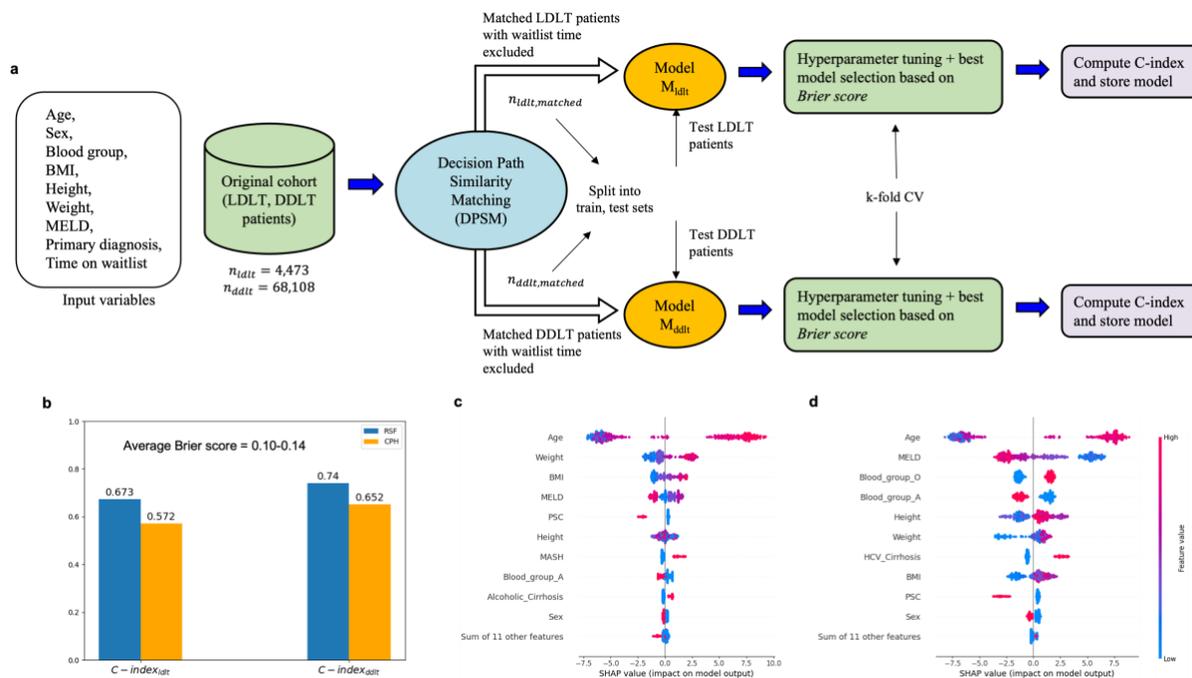


Fig. 3. Performance evaluation and SHAP.

(a) Model training and evaluation methodology: LDLT, DDLT cohorts are passed into the DPSM matching algorithm using all selected variables to account for confounding. Matched cohorts with waitlist time excluded are split into train and test sets. 2 graft-specific models are trained independently on the training samples, using a k-fold ($k=5$) cross-validation strategy and the best model across various hyperparameter settings is selected based on minimum Brier score. C-index is computed on the validation set and the best model is saved for further evaluation on the held-out test set; (b) RSF model C-index computed on the test set. Random Survival Forest (RSF) performance is compared with Cox proportional hazards (CPH) model, used as a baseline. Average Brier score varied from 0.10-0.14; top 10 features as predicted by SHAP on test set patients; (c) RSF-LDLT model applied to LDLT patients; (d) RSF-DDLT model applied to matched DDLT patients.

182

183 2.4. Estimation of LDLT benefit

184 A major strength of our ML-based technique is its ability to personalize survival prediction by
 185 incorporating individual patient variables into the modeling process. The emphasis in this
 186 section is on estimating the benefit in receiving an LDLT over a DDLT for an individual patient.
 187 To achieve this, we applied the trained RSF-LDLT and RSF-DDLT models to the LDLT and their
 188 matched DDLT counterparts, respectively. By contrasting the survival predictions generated by
 189 these models, we can estimate the personalized benefit of LDLT for each patient.

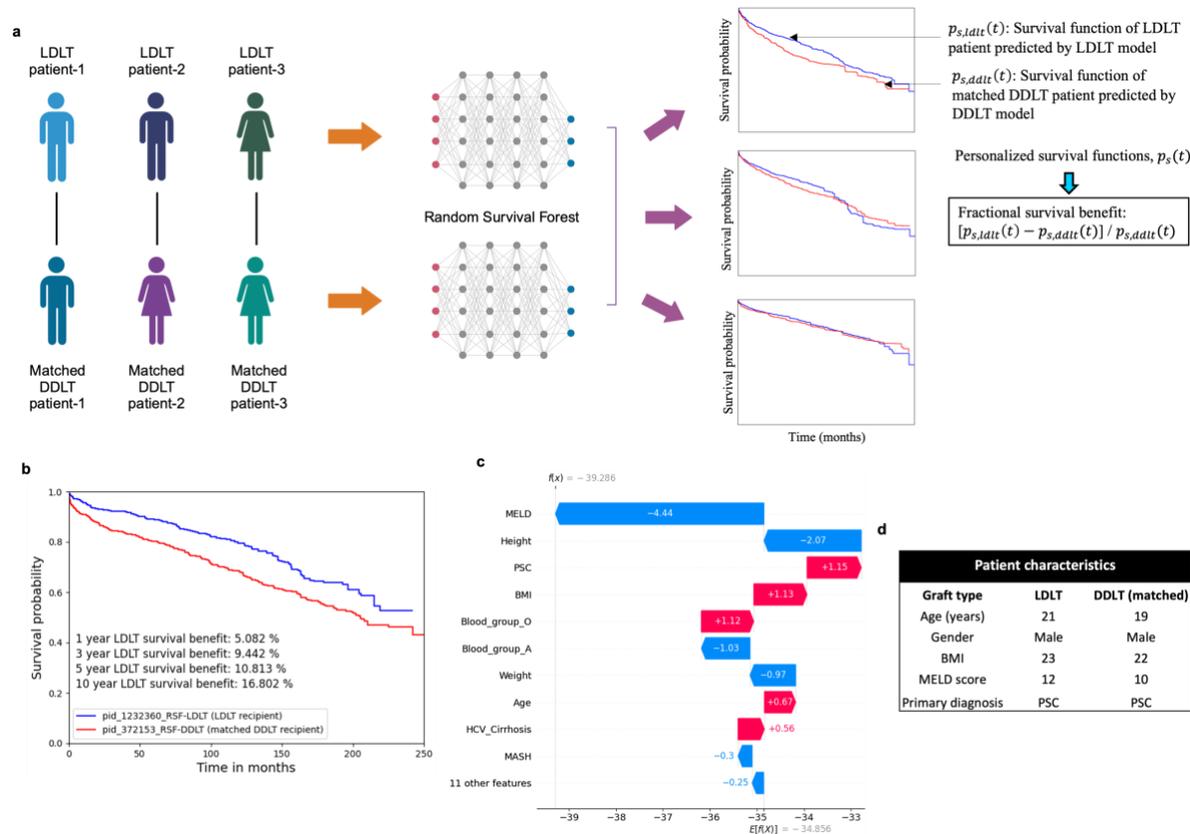


Fig. 4. LDLT versus DDLT survival for individual patient.

(a) Estimation of LDLT survival benefit: LDLT, DDLT survival models are applied independently to obtain the predicted survival functions $p_{s,ldlt}(t)$ and $p_{s,ddlt}(t)$ for LDLT and matched DDLT recipients (test), respectively. At a desired evaluation time-point, fractional LDLT survival benefit is defined as the probability difference between LDLT and DDLT survival normalized by DDLT survival; (b) LDLT benefit (%) at 1-, 3-, 5- and 10-years post-transplant for a patient diagnosed with PSC. Blue and red curves show survival functions of the LDLT and matched DDLT patient predicted by the respective models. Individual SHAP explanations (c) as well as corresponding patient characteristics (d) are also shown.

190
 191 Fig. 4 shows individual examples for two etiologies of interest: for the patient diagnosed with
 192 Primary Sclerosing Cholangitis (PSC), our model predicts a 10-year survival gain of 16.8 % (Fig.
 193 4b). To further understand the factors driving this predicted benefit, we conducted a patient-
 194 specific differential SHAP analysis (Fig. 4c), which identifies the key features contributing to
 195 survival gain. Variables highlighted in blue are associated with increased survival gain, while
 196 those in red indicate increased risk. The top variables that strongly influenced benefit were
 197 MELD, height, PSC and BMI. Original patient variables are also shown (Fig. 4d).

198
 199 Next, we compute the population-level predicted LDLT benefit by considering the matched
 200 patient groups. This is done on the held-out test set of patients, untouched during the model
 201 training process (Fig. 5a). Based on variables at the time of listing, our ML model predicts a

202 mean long-term (10-year) benefit of 10.3%. For validation, we also performed standard Kaplan-
 203 Meier analysis on the same set of patients to evaluate observed LDLT versus DDLT survival
 204 differences (Fig. 5b).

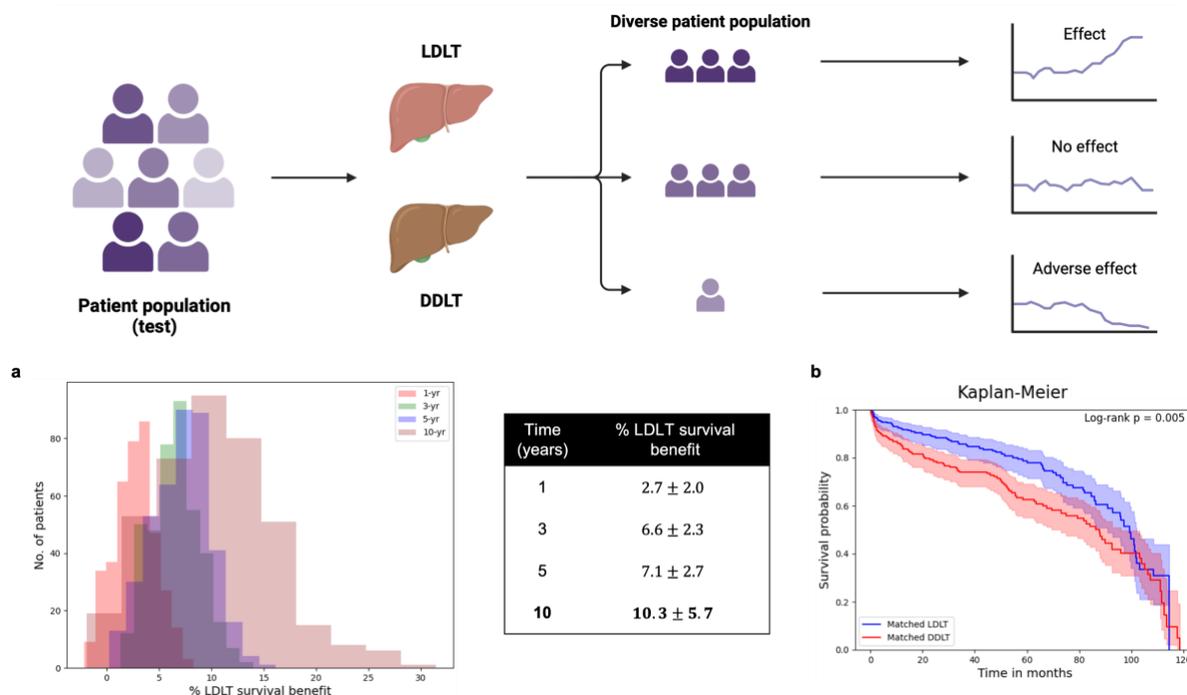


Fig. 5. LDLT survival benefit.

(a) 1, 3, 5 and 10-year fractional survival benefit of receiving LDLT over DDLT ($n_{test} = 401$) for LDLT recipients. If t does not exist for either of the LDLT, DDLT models, we perform interpolation, so that % benefit can be computed appropriately; (b) Kaplan-Meier estimator applied on the test set confirms LDLT survival benefit at the population-level. Survival times are cut-off at 120 months (10-years).

205
 206 **2.5. Etiology-specific benefit**

207 Finally, we evaluated the predicted survival benefit of LDLT over DDLT across six different
 208 etiologies (primary diagnoses) within the matched cohort, namely Autoimmune Hepatitis (AH),
 209 PBC, PSC, HCV Cirrhosis (HCV), MASH, and Alcoholic Cirrhosis (AC). By analyzing the
 210 survival outcomes for specific diseases, we aimed to identify which etiologies were associated
 211 with the highest LDLT benefit.

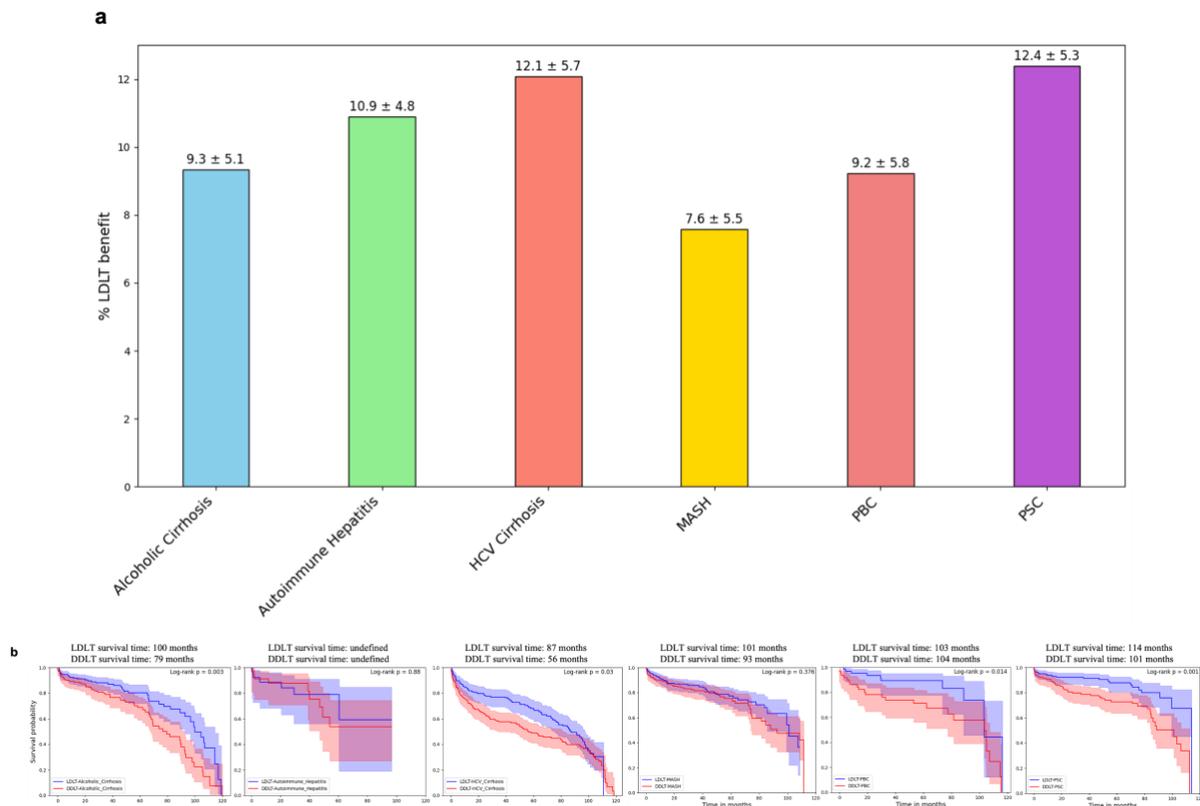


Fig. 6. Etiology-specific benefit.

(a) Bar chart shows 10-year (long-term) post-transplant LDLT benefit (%) for each of the 6 distinct etiologies: Alcoholic Cirrhosis, Autoimmune Hepatitis, HCV Cirrhosis, MASH, PBC and PSC. Results are computed on the held-out test set, these patients are untouched during RSF model training; (b) observed survival (Kaplan-Meier) for all 6 etiologies to evaluate matched LDLT (blue) vs DDLT (red) survival differences. Survival data is cut-off at 10 years.

212
 213 Using our ML-guided approach, we computed the differential survival gain of LDLT for each
 214 patient and then aggregated these results based on their primary diagnosis (Fig. 6a). The
 215 analysis revealed that certain etiologies exhibited a significantly higher survival benefit when
 216 transplanted with LDLT compared to DDLT. Patients diagnosed with PSC (12.4 ± 5.3 %) and
 217 HCV (12.1 ± 5.7 %) showed substantial long-term survival advantages with LDLT, over a 10-
 218 year period. For comparison with ground truth, we also evaluated observed survival differences
 219 between the two groups (Fig. 6b). These findings suggest that LDLT may be particularly
 220 advantageous for patients with these conditions, potentially influencing clinical decision-
 221 making.

222

223 **3. Discussion**

224 The challenge of determining which patients would benefit most from receiving a living donor
225 liver transplant (LDLT) versus a deceased donor liver transplant (DDLT) is compounded by the
226 impracticality of conducting randomized controlled trials (RCTs). To address this, we create a
227 novel approach which we call Decision Path Similarity Matching (DPSM) algorithm, that
228 combines advanced matching and machine learning (ML) techniques to emulate the balanced
229 conditions of an RCT using observational data to personalize survival predictions. By effectively
230 minimizing confounding factors, DPSM enables causal-type estimation of the effect of graft-type
231 on post-transplant survival.

232

233 **3.1. DPSM: a novel, multivariate method for one-to-one matching**

234 Our innovative algorithm represents a significant advancement in the field of observational study
235 design, particularly in contexts where conducting randomized controlled trials (RCTs) is
236 impractical or impossible. Unlike traditional Propensity Score Matching (PSM), which
237 compresses complex patient data into a single propensity score, DPSM leverages the full
238 decision paths generated by a Random Forest classifier to perform one-to-one matching
239 between LDLT and DDLT patients. This approach allows DPSM to retain the multidimensional
240 complexity of patient data, resulting in more nuanced and accurate matching that closely mimics
241 the balance achieved in RCTs.

242

243 The key strength of DPSM is its ability to minimize confounder bias, thereby producing more
244 balanced cohorts than traditional PSM. Our study demonstrated that DPSM significantly
245 decreased the AUROC for graft-type prediction after matching, a clear indication that DPSM
246 more effectively mitigates covariate differences between the LDLT and DDLT groups. This
247 enhanced performance is crucial in creating a robust foundation for subsequent survival
248 analysis, ensuring that any observed differences in outcomes are attributable to graft-type
249 differences.

250

251 Additionally, by utilizing decision paths instead of a single scalar score, DPSM allows clinicians
252 and researchers to understand the prominent variables involved in the decision-making process
253 leading up to the matching. These insights are particularly valuable in clinical settings, where
254 understanding the rationale behind matching decisions can foster greater confidence in the
255 study's findings and support more informed clinical decision-making. The ability to visualize and
256 interpret the decision-making process also aligns DPSM more closely with the principles of
257 RCTs, where the reasoning behind patient assignment is clear and systematic.

258

259 In summary, DPSM is a powerful technique, enabling researchers to simulate the conditions of
260 an RCT more effectively in observational studies. Its superior matching performance and
261 transparency make it a valuable tool not only for advancing research in liver transplantation but
262 also for broader applications where the target questions are causal in nature. As the field of
263 clinical research increasingly turns to observational data in the absence of feasible RCTs,
264 methods like DPSM will play a critical role in ensuring that the insights drawn from these studies
265 are both accurate and actionable.

266

267 **3.2. ML framework for personalized survival predictions**

268 Applying our RSF model on a test set of patients, we report a mean predicted LDLT survival
269 gain of 2.7%, 6.6%, 7.1% and 10.3% at 1-, 3-, 5- and 10-years post-transplant respectively. It is
270 worth placing this in the context of prior evidence pertaining to LDLT versus DDLT survival.
271 Barbetta *et. al.*¹⁰ also analyzed SRTR data and found a mortality risk reduction of 17%, 15%
272 and 13% at 1-, 3- and 5-year post-transplant for LDLT recipients. Higher estimated benefit in the
273 latter is suspected to be a consequence of the lack of matching in their analysis. In fact, another
274 study¹¹ around the same time that analyzed Canadian transplant recipients found that graft-type
275 differences got washed away upon adjustment for donor as well as recipient characteristics.

276

277 In our study, where we predict individual LDLT benefit using patient-specific variables at the
278 time of listing, we find significant heterogeneity across all patients that we evaluated our method
279 on, underscoring the importance of personalization. The ability to predict individualized survival
280 outcomes allows clinicians to move beyond a 'one-size-fits-all' approach, enabling more tailored
281 decisions that align with the specific characteristics and needs of each patient. This not only
282 optimizes transplant outcomes but also enhances patient counseling, as clinicians can provide
283 more accurate, data-driven information to patients and their families when discussing treatment
284 options.

285

286 Finally, our etiology-specific analysis underscores the importance of individualized treatment
287 planning, as the survival benefit of LDLT can vary significantly depending on the underlying liver
288 disease. These insights could guide clinicians in making more informed decisions about
289 transplant strategies, particularly for patients with specific diagnoses where LDLT offers the
290 most substantial benefit. Among all etiologies tested, our ML tool predicted the highest benefit
291 for patients diagnosed with PSC (12.4 ± 5.3 %). These patients often experience slow

292 progression of the disease, but when complications such as cholangitis arise, timely
293 transplantation becomes critical. LDLT, with its shorter wait times, offers a significant survival
294 advantage in such urgent cases. Our result is backed by a recent study by Sierra *et. al.*¹⁷, which
295 also reported a long-term (10-year) survival advantage for PSC patients who received an LDLT
296 (81.9 %) over a DDLT (72.7 %).

297
298 Integrating these personalized predictions into clinical workflows could significantly improve
299 decision-making processes, ensuring that patients are referred for LDLT especially when the
300 survival benefit is significant. As the field of liver transplantation progresses towards precision-
301 based approaches, tools like these will become essential in guiding downstream clinical
302 decisions and potentially improving overall patient care.

303

304 **3.3. Study strengths and limitations**

305 This study has several notable strengths that enhance its contribution to the field of liver
306 transplantation. First, the development and application of the Decision Path Similarity Matching
307 (DPSM) algorithm represents a significant advancement in observational study design, allowing
308 us to create well-matched cohorts that closely mimic the conditions of an RCT. This is evident in
309 the substantial reduction in systematic differences between the LDLT and DDLT cohorts, as
310 evidenced by the drop in AUROC from 0.83 (pre-matching) to 0.51 (post-matching), which
311 shows that our algorithm effectively removed confounding associations among variables. By
312 reducing bias and aligning survival outcomes across groups, DPSM enables a more accurate
313 comparison of LDLT and DDLT outcomes, which is crucial in the absence of feasible RCTs.
314 Additionally, the integration of Random Survival Forest (RSF) models to provide personalized
315 survival predictions adds a valuable dimension to clinical decision-making, offering tailored
316 insights that can optimize patient care.

317

318 Our study also has important limitations. The retrospective nature of the analysis, relying on
319 data from the SRTR, may introduce biases inherent in observational studies. It is important to
320 note that DPSM is able to mitigate these confounding biases, as long as they are observable,
321 i.e., captured within our dataset. Another limitation is the reliance on clinico-demographic
322 variables available at the time of listing, which, although comprehensive, may not capture all
323 factors influencing transplant outcomes. Incorporating additional variables, such as genetic
324 markers or more detailed comorbidity data, could further refine the predictive models and
325 enhance the precision of survival estimates.

326
327 We also acknowledge the absence of external validation as a limitation of this study. This was
328 primarily due to the lack of access to sufficiently large datasets outside the SRTR. DPSM
329 performs optimally with larger sample sizes, especially the DDLT pool, as these allow for more
330 reliable transformation of observational data into RCT-like conditions. Smaller datasets may not
331 provide the robustness needed for effective matching, underscoring the importance of future
332 research focused on validating these findings in different cohorts and clinical environments.
333 Finally, while the RSF model demonstrated good performance in this study, further validation in
334 prospective studies and across different cohorts is necessary to confirm its broader applicability.

335
336 In conclusion, this study introduces a novel Decision Path Similarity Matching (DPSM)
337 methodology, which represents a significant advancement in creating RCT-like comparability
338 from observational data and debiasing transplant outcomes. DPSM offers a more transparent
339 and explainable approach to matching patients compared to traditional methods, allowing for
340 personalized predictions of survival benefit in living donor liver transplantation. While further
341 research and external validation are required to enhance its robustness and generalizability, the
342 innovations presented here mark an important step toward more individualized, data-driven
343 decision-making in liver transplantation.

344

345 **4. Methods**

346 **4.1. Study Design**

347 We conducted a retrospective, cross-sectional study from the SRTR database. The SRTR
348 includes data of all organ donors, as well as waitlisted and transplanted recipients in the United
349 States submitted by the Organ Procurement and Transplantation Network (OPTN). This study
350 was approved by the Research Ethics Board at the University Health Network.

351
352 Adult (≥ 18 -years-old) LT patients, listed between 28th February 2002 and 23rd May 2023, were
353 included in the analysis. Study exclusion criteria is clearly defined in Fig. S1. Patients with
354 reported MELD scores greater than 40 were excluded. These patients (<5% of the entire
355 population) are generally very sick and highly prone to pre-operative mortality, including them
356 would make it challenging to clearly delineate the effect of transplant type on post-transplant
357 survival. Recipients with previous or multi-organ transplants were excluded as were those

358 diagnosed with HIV, acute liver failure (ALF) and HCC. Patients who received exception points
359 were also removed.

360

361 **4.2. Variables and Outcomes**

362 Clinico-demographic patient variables: age, sex, blood group, BMI, height, weight, MELD score,
363 primary diagnosis (indication for transplant) and time on the waitlist were collected at the time of
364 listing. For variable comparison between the LDLT, DDLT patient cohorts, an alpha level of
365 <0.05 was selected as the significance threshold. Post-transplant, mortality (all-cause) was the
366 outcome or event of interest and event times were defined from the time of transplant to either
367 the time of death or censored at the date of last follow-up.

368

369 **4.3. Data Preparation and Preprocessing**

370 Covariates with greater than 20% missingness were excluded from the analysis. For continuous
371 variables with missing values, mean imputation was performed. Categorical variables were one-
372 hot encoded (OHE). Subsequently, we ended up with 21, unit normalized input features to the
373 ML model.

374

375 **4.4. Matching**

376 We designed and implemented a new method, Decision Path Similarity Matching (DPSM) to
377 account for confounding bias, ensuring that predicted survival differences would be
378 predominantly attributed to graft-type. The key feature of DPSM being that it uses the “similarity”
379 or closeness between decision paths, which provides a richer feature encoding as opposed to
380 matching based on output probabilities alone in the case of PSM. The constituent steps of our
381 algorithm are shown in Fig. 1a along with the pseudocode below (Fig. 7) - (1) first, the original
382 LDLT-DDLT dataset is randomly split into train (70%) and test (30%) sets. A Random Forest
383 (RF) classifier is trained on the training set to predict transplant type using all input variables
384 previously listed in Sec. 4.2. RF was selected due to its ability to capture non-linear
385 relationships among variables and handle imbalanced data, crucial in the LDLT versus DDLT
386 context. (2) The best model is selected across a hyperparameter search ($n_{estimators} =$
387 $50, 100, 200, 300, 500, 1000$, $min_{samples, leaf} = 500$) using $k(=5)$ -fold cross-validation (cv) on the
388 train set and performance is evaluated on the test set using the area under the receiver operator
389 characteristic (AUROC) curve. (3) From the trained model, we extract decision paths for
390 individual patients in the held-out test set, averaged across all the Trees in the Forest. This is an

391 n-dimensional binary vector [1, 0, 1, 0, ..], where “n” is the total number of decision nodes
392 (features) per tree. A node is “1” if it applies to a particular patient and “0” otherwise. (4)
393 Hamming distance (d_H) is then computed for every pair of (LDLT, DDLT) decision paths. To
394 remove “poor” matches or outliers, patients whose pairwise d_H exceeds a selected threshold
395 (th) are removed. The optimal threshold was determined as that which minimized
396 $|AUROC_{post-match} - 0.5|$, at an acceptable patient dropout rate (Table S1). We selected $th =$
397 0.45. (5) The filtered distance matrix is used to perform one-to-one matching using the “Munkres
398 Assignment” procedure¹⁸. This is an “optimal” algorithm that matches by minimizing the total
399 Hamming distance across all LDLT-DDLT patient pairs. After this step, we repeat step (2) to
400 calculate the AUROC for the matched dataset using a k(=5)-fold cv strategy. The same set of
401 hyperparameters are used. Finally, match effectiveness is quantified as the difference between
402 pre- and post-match AUROCs.

Input:

LDLT patients' data (X_{ldlt})

DDLT patients' data (X_{ddlt})

Trained Random Forest model (M_{rf}) to predict graft-type

1. For each tree t in RandomForest M_{rf} :
 - a. Extract decision paths for LDLT, DDLT patients from tree
 $Paths_{ldlt}[t] = \text{decision_path}(X_{ldlt}, t)$
 $Paths_{ddlt}[t] = \text{decision_path}(X_{ddlt}, t)$
2. For each tree t in RandomForest M_{rf} :
 - a. Compute vectorized Hamming distances for each LDLT-DDLT pair:
 $Distance[t] = \text{hamming_distance}(Paths_{ldlt}[t], Paths_{ddlt}[t])$
3. Compute AvgDistance for all patient pairs by averaging across all trees:
 $AvgDistance = \text{average}(Distance[t] \text{ for all } t)$
4. Filter patients with bad matches using optimized threshold (th)
 $FilteredAvgDistance = AvgDistance[AvgDistance > th]$
5. Perform optimal matching between LDLT and DDLT patients using Munkres Assignment procedure:
 $MatchedPairs = \text{linear_sum_assignment}(AvgDistance)$
6. Evaluate matching & visualize results.

Output: Matched LDLT-DDLT patient pairs

Fig. 7. DPSM: pseudocode

403

404

405 4.5. Model Training and Validation

406 The matched sub-populations ($n_{ldlt} = n_{ddlt} = 1,337$) were first split into train (70%) and test
407 (30%) sets ($n_{train} = 936, n_{test} = 401$). For the time-to-event survival prediction, two Random
408 Survival Forest (RSF) models were trained independently on the matched LDLT-DDLT cohorts,
409 respectively. Hyperparameter tuning was performed on the train dataset using k (=5)-fold cross
410 validation and the model that produced minimum Brier score was selected as the best one. For
411 validation, both C-index and Brier score, averaged across 5 evaluation time points: 0.5, 1, 3, 5
412 and 10 years were computed on the held-out test set.

413 During hyperparameter tuning, it is common to use C-index as the evaluation metric for time-to-
414 event models. While this is useful in understanding relative risk ranking, it is not informative
415 about the accuracy and calibration of the predicted survival predictions. We instead utilize the
416 Brier score, defined as follows: $BS = (1/N) \sum_{i=1}^N (gt_i - p_i)^2$, where N is the sample size, gt_i and
417 p_i are the actual and predicted event probabilities for observation i . This metric quantifies the
418 error rate between prediction and ground truth, serving as an ideal choice for calibrating survival
419 models. Subsequently, the optimal hyperparameters selected are those that minimize the time-
420 averaged Brier score.

421

422 **4.6. Estimation of LDLT Survival Benefit**

423 For a given LDLT and matched DDLT patient, the corresponding trained models (RSF-LDLT,
424 RSF-DDLT) are applied to obtain the respective predicted survival functions $S_{ldlt}(t)$ and
425 $S_{ddlt}(t)$. RSF models produce unique event times according to the data they were trained on.
426 To ensure that both LDLT and DDLT models predict definitive survival for a given evaluation
427 time point t_i , we interpolate predicted survival probabilities across all time. Finally, differential
428 LDLT benefit is evaluated as $(S_{ldlt}(t) - S_{matched\ ddt}(t))/S_{matched\ ddt}(t)$.

429

430 **5. Code availability**

431 The source code for this work is available on GitHub
432 (https://github.com/Anivader/LDLT_survival_benefit_ML_tool). All analysis was performed using
433 Python.

434

435 **6. References**

- 436 1. Spieth, P. M. *et al.* Randomized controlled trials – a matter of design. *Neuropsychiatr.*
437 *Dis. Treat.* **12**, 1341–1349 (2016).
- 438 2. Schiano, T. D., Kim-Schluger, L., Gondolesi, G. & Miller, C. M. Adult living donor liver
439 transplantation: The hepatologist’s perspective. *Hepatology* **33**, 3–9 (2001).
- 440 3. Brown, J., Sorrell, J. H., McClaren, J. & Creswell, J. W. Waiting for a Liver Transplant.
441 *Qual. Health Res.* **16**, 119–136 (2006).
- 442 4. Larson, A. M. & Curtis, J. R. Integrating Palliative Care for Liver Transplant
443 Candidates“Too Well for Transplant, Too Sick for Life”. *JAMA* **295**, 2168–2176 (2006).

- 444 5. Humar, A. *et al.* Adult Living Donor Versus Deceased Donor Liver Transplant (LDLT
445 Versus DDLT) at a Single Center: Time to Change Our Paradigm for Liver Transplant. *Ann.*
446 *Surg.* **270**, 444 (2019).
- 447 6. Ivanics, T. *et al.* Low utilization of adult-to-adult LDLT in Western countries despite
448 excellent outcomes: International multicenter analysis of the US, the UK, and Canada. *J.*
449 *Hepatol.* **77**, 1607–1618 (2022).
- 450 7. Tran, L. & Humar, A. Expanding living donor liver transplantation in the Western world:
451 changing the paradigm. *Dig. Med. Res.* **3**, (2020).
- 452 8. Karnam, R. S. *et al.* Sex Disparity in Liver Transplant and Access to Living Donation.
453 *JAMA Surg.* **156**, 1010–1017 (2021).
- 454 9. Lee, S.-G. A Complete Treatment of Adult Living Donor Liver Transplantation: A Review
455 of Surgical Technique and Current Challenges to Expand Indication of Patients. *Am. J.*
456 *Transplant.* **15**, 17–38 (2015).
- 457 10. Barbeta, A. *et al.* Meta-analysis and meta-regression of outcomes for adult living donor
458 liver transplantation versus deceased donor liver transplantation. *Am. J. Transplant. Off. J.*
459 *Am. Soc. Transplant. Am. Soc. Transpl. Surg.* **21**, 2399–2412 (2021).
- 460 11. Goto, T. *et al.* Superior Long-Term Outcomes of Adult Living Donor Liver
461 Transplantation: A Cumulative Single-Center Cohort Study With 20 Years of Follow-Up. *Liver*
462 *Transpl.* **28**, 834–842 (2022).
- 463 12. Benedetto, U., Head, S. J., Angelini, G. D. & Blackstone, E. H. Statistical primer:
464 propensity score matching and its alternatives†. *Eur. J. Cardiothorac. Surg.* **53**, 1112–1117
465 (2018).
- 466 13. Abadie, A. & Imbens, G. W. Matching on the Estimated Propensity Score. *Econometrica*
467 **84**, 781–807 (2016).
- 468 14. King, G. & Nielsen, R. Why Propensity Scores Should Not Be Used for Matching. *Polit.*
469 *Anal.* **27**, 435–454 (2019).
- 470 15. Nguyen, T.-L. *et al.* Double-adjustment in propensity score matching analysis: choosing
471 a threshold for considering residual imbalance. *BMC Med. Res. Methodol.* **17**, 78 (2017).
- 472 16. Rubin, D. B. & Thomas, N. Combining Propensity Score Matching with Additional
473 Adjustments for Prognostic Covariates. *J. Am. Stat. Assoc.* **95**, 573–585 (2000).
- 474 17. Sierra, L. *et al.* Living-Donor Liver Transplant and Improved Post-Transplant Survival in
475 Patients with Primary Sclerosing Cholangitis. *J. Clin. Med.* **12**, 2807 (2023).
- 476 18. Munkres, J. Algorithms for the Assignment and Transportation Problems. *J. Soc. Ind.*
477 *Appl. Math.* **5**, 32–38 (1957).

478

479 **7. Acknowledgements**

480 Mamatha Bhat acknowledges support from the Toronto General and Western Hospital
481 Foundation, Canadian Institutes for Health Research and Canadian Donation and Transplant
482 Research Program. Chris McIntosh holds the Chair in Medical Imaging at the Joint Department
483 of Medical Imaging at the University Health Network, and the Department of Medical Imaging at
484 the University of Toronto. Michael Brudno holds a CIFAR AI Chair.

485

486 **8. Author contributions**

487 A.G., C.M., M. Bhat. and Y.S. conceptualized the study. C.M. supervised the experimental
488 design, A.G. developed the computational analysis pipelines and generated all the data. Y.S.
489 helped with the data pre-processing script. B.J.H. and X.Z. helped with clinical interpretability.
490 A.G. wrote the manuscript and C.M., B.J.H., M. Bhat provided feedback and M. Brudno
491 reviewed the manuscript. C.M. and M. Bhat. supervised the study.

492

493 **9. Ethics declarations**

494 The authors declare no competing interests.