

**An AI system for the detection and incidence prediction of chronic kidney disease and type 2 diabetes based on retinal fundus images**

Kang Zhang<sup>1\*#</sup>, Xiaohong Liu<sup>2\*</sup>, Jin Yuan<sup>3\*</sup>, Jie Xu<sup>4\*</sup>, Wenjia Cai<sup>3\*</sup>, Kai Wang<sup>5</sup>, Ting Chen<sup>2#</sup>, Yuanxu Gao<sup>1</sup>, Sheng Nie<sup>6</sup>, Xiaoqi Qin<sup>5</sup>, Wenqin Xu<sup>1</sup>, Andrea Olvera<sup>1</sup>, Kanmin Xue<sup>7</sup>, Zhihuan Li<sup>1</sup>, Yuandong Su<sup>8</sup>, Meixia Zhang<sup>8</sup>, Charlotte L. Zhang<sup>9</sup>, Oulan Li<sup>9</sup>, Edward E. Zhang<sup>9</sup>, Jie Zhu<sup>10</sup>, Yiming Xu<sup>2</sup>, Daniel Kermany<sup>1</sup>, Kaixin Zhou<sup>9</sup>, Ying Pan<sup>10</sup>, Shaoyun Li<sup>12</sup>, Iat Fan Lai<sup>13</sup>, Ying Chi<sup>14</sup>, Changuang Wang<sup>15</sup>, Qi Zhang<sup>16</sup>, Johnson Lau<sup>17</sup>, Dennis Lam<sup>18</sup>, Yin Shen<sup>19</sup>, Tao Xu<sup>9#</sup>, Yong Zhou<sup>20</sup>, and Guangyu Wang<sup>5#</sup>

<sup>1</sup> Center for Biomedicine and Innovations, Faculty of Medicine, Macau University of Science and Technology, Macau, China

<sup>2</sup> Department of Computer Science and Technology, Tsinghua University, Beijing, China

<sup>3</sup> State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou, China

<sup>4</sup> Beijing Institute of Ophthalmology, Beijing Tongren Eye Center, Beijing Tongren Hospital, Beijing Ophthalmology and Visual Science Key Lab, Beijing, China

<sup>5</sup> School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing, China.

<sup>6</sup> State Key Laboratory of Organ Failure Research; National Clinical Research Center for Kidney Disease & Nanfang hospital, Southern Medical University, Guangzhou, China

<sup>7</sup> Nuffield Department of Neuroscience, Oxford University, Oxford, UK

<sup>8</sup> Center for Translational Innovations, West China Hospital and Sichuan University, Chengdu, China

<sup>9</sup> Bioland Laboratory, Guangzhou, China

<sup>10</sup> Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, China

<sup>11</sup> Department of Endocrinology, Kunshan Hospital Affiliated to Jiangsu University, Kunshan, China

<sup>12</sup> The Big Data Research Center, Chongqing Renji affiliated Hospital to the University of Chinese Academy of Sciences, Chongqing, China.

<sup>13</sup> Ophthalmic Center, Kiang Wood Hospital, Macau, China

<sup>14</sup> Peking University First Affiliated Hospital, Beijing, China

<sup>15</sup> Peking University Third Affiliated Hospital, Beijing, China

<sup>16</sup> Biotherapy Center, Third Affiliated Hospital of Sun Yat-sen University, Guangzhou, China

<sup>17</sup> Department of Applied Biology and Chemical Technology, Hong Kong Polytechnic University, Hong Kong, China

<sup>18</sup> C-MER Dennis Lam and Partners Eye Center, C-MER International Eye Care Group, Hong Kong, China

<sup>19</sup> Medical Research Institute, Wuhan University, Wuhan, China

<sup>20</sup> Clinical Research Center, Shanghai First People's Hospital, Shanghai Jiaotong University, Shanghai, China

# These authors contributed equally to this work.

\*Correspondence and reprints should be addressed to: [guangyu.wang24@gmail.com](mailto:guangyu.wang24@gmail.com); [tingchen@mail.tsinghua.edu.cn](mailto:tingchen@mail.tsinghua.edu.cn); [xutao@ibp.ac.cn](mailto:xutao@ibp.ac.cn); or [kang.zhang@gmail.com](mailto:kang.zhang@gmail.com)

Editorial corresponding author:

49 Kang Zhang, MD, PhD  
50 Center for Biomedicine and Innovations, Faculty of Medicine  
51 Macau University of Science and Technology  
52 Email address: [kang.zhang@gmail.com](mailto:kang.zhang@gmail.com)  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65  
66  
67  
68  
69  
70  
71  
72  
73  
74  
75  
76  
77  
78  
79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98  
99  
100  
101

102     **Abstract**

103

104     Regular health screening plays a crucial role in the early detection of common chronic  
105     diseases and prevention of their progression. An AI system capable of recapitulating  
106     early disease detection, staging and incidence prediction would help to improve  
107     healthcare access and delivery, particularly in resource-poor or remote settings. Using  
108     a total of 115,344 retinal fundus photographs from 57,672 patients (with data split  
109     into mutually exclusive training, internal testing, and external validation sets), we first  
110     developed AI models capable of identifying chronic kidney disease (CKD) and type 2  
111     diabetes mellitus (T2DM) based on fundus images. The AI system was shown to be  
112     capable of predicting the clinical indicators of CKD and T2DM (including eGFR and  
113     blood glucose levels), which indicates its potential for extracting quantitative clinical  
114     metrics embedded subtly within retinal fundus images. We further developed an AI  
115     system to predict the risk of disease progression using baseline images of 10,269  
116     patients for whom longitudinal clinical data were available for up to 6 years, which  
117     demonstrated potential utility in optimizing health screening intervals and clinical  
118     management. The generalizability of the AI system in identifying and predicting the  
119     progression of CKD and T2DM was evaluated using population-based external  
120     validation cohorts. Moreover, a prospective pilot study with 3,081 patients was also  
121     conducted to demonstrate the broader applicability of the AI system at the  
122     ‘point-of-care’ using fundus images captured with smartphones. The results provide  
123     proof-of-concept for a reliable and non-invasive AI-based clinical screening tool  
124     based on fundus photographs for the early detection and incidence prediction of two  
125     common systemic diseases.

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

## Introduction

Systemic diseases, including chronic kidney disease (CKD) and diabetes mellitus, pose major healthcare challenges. CKD is a highly prevalent disease and affects approximately 8-16% of the world population<sup>1,2</sup>. CKD is a serious public health problem, as its adverse outcome is not limited to end-stage renal failure requiring dialysis or transplantation but also vascular complications of impaired kidney function<sup>3</sup>. Moreover, cardiovascular events and mortality are strongly associated with CKD in the high-risk diabetic or hypertensive population. Type 2 diabetes mellitus (T2DM) is another major common chronic disease globally, with an estimated prevalence of 9.3% (463 million affected individuals) in 2019. According to the International Diabetes Federation, its prevalence has been increasing steadily in recent years and will reach an estimated 700 million by 2045<sup>4</sup>. According to the US Center for Disease Control and Prevention, diabetes is one of the leading causes of mortality globally. It is also a leading risk factor for many other common medical problems, including cardiovascular disease, kidney failure and blindness<sup>5-7</sup>. In many of these conditions, early diagnosis and treatment are crucial in reducing the associated comorbidities and mortality. However, early identification and diagnosis of CKD and diabetes remain challenging as many patients are asymptomatic or only have non-specific symptoms, with some reports suggesting that up to 5% of the diabetic population remains undiagnosed.

The American Diabetes Association (ADA)<sup>8</sup> recommends diabetic patients with over 5-year disease course to an annual test on urinary albumin and estimated glomerular filtration rate (eGFR). A previous study reported a less than 20% awareness rate of CKD and a less than 50% treatment rate<sup>9</sup> in China. Regular screening is crucial for the early detection and diagnosis of CKD and the prevention of its progression. A cost-benefit analysis conducted in China revealed that medical expenses could be considerably reduced by carrying out kidney disease screening in newly diagnosed T2DM patients<sup>10</sup>. Early detection and intervention are key to the prevention of end-stage renal failure in CKD and sight-threatening complications of diabetic retinopathy.

The retina of the eye is a convenient window into homeostasis of the body, in which we can non-invasively observe vascular, neural, and connective tissues, both structurally and, in the case of the vasculature, in dynamic action. Systemic disorders may have manifestations in the fundus that allow us to detect, diagnose, stage, monitor, and manage the systemic disease. Recent advances using deep learning classifiers have led to applications of artificial intelligence (AI) in many areas of healthcare<sup>11-15</sup>, including image-based diagnosis<sup>16</sup> and natural language processing<sup>17</sup>. In particular, convolutional neural networks with transfer learning have facilitated efficient and accurate image-based diagnosis well beyond human capabilities<sup>16,18</sup>. There has also been early but promising evidence that identification of systemic conditions and clinical metrics based on fundus photographs is possible<sup>19,20</sup>,

190 suggesting that deep learning algorithms can detect subtle associations that are  
191 undetectable to the human observer. Therefore, retinal images, which can be acquired  
192 rapidly and non-invasively, may have the potential to provide ‘point-of-care’  
193 biomarkers for systemic disease.

194  
195 Though there is a recent report on advanced CKD detection using a retinal  
196 image-based AI system<sup>21</sup>, prediction of CKD onset from a normal baseline and  
197 diagnosis of early CKD (proteinuria) have not yet been described and would have an  
198 important impact on disease prevention and outcomes. Similarly, predicting the onset  
199 of type 2 diabetes would be critical for disease prevention and improving outcomes.  
200 Here, we explored risk stratification for developing CKD and T2DM using both  
201 retinal images and known clinical risk factors. Identifying this asymptomatic  
202 pre-morbid population provides the novel possibility of better channeling healthcare  
203 resources to monitor ‘patients-at-risk’ and to modify lifestyle and other risk factors at  
204 an early stage.

205  
206 One of the key practical challenges to the application of AI in healthcare, particularly  
207 in low-resource settings, is the lack of stable computational infrastructure and  
208 resources required to run the AI algorithms. Accordingly, the deployment of AI-based  
209 technologies through mobile platforms has emerged as a growing area of  
210 investigation<sup>22</sup>. The latest personal smartphones are equipped with the requisite  
211 hardware to run AI algorithms, such as Google Translate, Siri, FaceID and shopping  
212 apps with object recognition. With Android’s Neural Network API (NNAPI) and the  
213 Neural Engine chip in the iPhone X (Apple Inc, Cupertino, CA, USA)<sup>23</sup>, smartphones  
214 now deploy a range of machine learning algorithms from image classification<sup>24</sup> to  
215 object tracking<sup>25</sup>, face detection and recognition<sup>26</sup>, natural language translation<sup>27</sup>,  
216 sentence completion<sup>28</sup>, and augmented reality<sup>29</sup>. Furthermore, smartphone ownership  
217 is widespread worldwide. For instance, in Nigeria, the physician-to-patient ratio is 1:  
218 2660<sup>30</sup>, while smartphone ownership is 1:3.5<sup>31</sup>. An AI-based smartphone diagnostic  
219 platform is thus an attractive way to broaden healthcare access by encouraging  
220 patients to self-monitor and allowing doctors to diagnose and follow-up patients  
221 remotely.

222  
223 In this study, we aimed to develop an AI system capable of analyzing retinal fundus  
224 images to detect CKD and T2DM. We employed deep learning-based analysis of  
225 retinal fundus images for two types of tasks: a regression task of predicting  
226 continuous values (including eGFR) and a binary classification task of making the  
227 diagnosis. We also validated the AI algorithm to detect CKD and T2DM in external  
228 independent patient populations (**Fig.1a**). Furthermore, we show our AI system can  
229 predict disease development and perform risk stratification for CKD and T2DM using  
230 retinal images in two longitudinal cohorts (**Fig. 1b**). Using this approach, targeted  
231 screening to subgroups of the population could potentially help deliver risk-reduction  
232 interventions to those most likely to develop the diseases. Finally, a smartphone-based

233 system was created to provide a ‘point-of-care’ platform for CKD and T2DM  
234 screening in the community.

235

236

237

## 238 **Results**

### 239 **Definitions of chronic kidney disease and type 2 diabetes**

240

241 Based on international guidelines and previous studies<sup>1,32</sup>, identifying individuals with  
242 CKD relies primarily on estimated glomerular filtration rate (eGFR), an index of  
243 kidney function, and renal damage markers (e.g. urinary albumin). The presence of  
244 CKD was defined by an eGFR  $\geq 60$  mL/min/1.73m<sup>2</sup> with albuminuria or eGFR  $< 60$   
245 mL/min/1.73 m<sup>2</sup>, confirmed in at least two visits separated by three months; whereas  
246 normal controls were defined by an eGFR  $\geq 60$  mL/min/1.73m<sup>2</sup> without albuminuria.  
247 In our study, we utilized the images and corresponding eGFR measurements of  
248 already-diagnosed CKD patients.

249

250 Furthermore, CKD can be divided into 5 stages depending on severity during clinical  
251 applications. For the purpose of this study, we defined three risk stratifications  
252 depending on the 5 stages: early CKD (stage 1 and 2, eGFR  $\geq 60$  mL/min/1.73m<sup>2</sup> with  
253 albuminuria); advanced CKD (stage 3, eGFR 30-59 mL/min/1.73m<sup>2</sup>), and severe+  
254 CDK (stage 4 and 5, eGFR  $< 30$  mL/min/1.73m<sup>2</sup>). Early identification of early CKD to  
255 advanced CKD is clinically important as a timely intervention to modify known risk  
256 factors (e.g. treatment of hypertension, diabetes and urinary tract obstructions) could  
257 slow disease progression and reduce long-term health and economic burdens.  
258 Additionally, one of the important goals of CKD management is to prevent end-stage  
259 kidney failure, which is associated with the need for renal dialysis or transplantation,  
260 and increased mortality<sup>33</sup>.

261

262 For the study of the other systemic disease of T2DM, we utilized the images and  
263 corresponding clinical data, including laboratory values of already-diagnosed T2DM  
264 patients, based on: fasting blood glucose level  $\geq 7.0$  mmol/L confirmed in at least two  
265 visits, or an HbA1c value  $\geq 6.5\%$ , and/or based on a history of drug treatment for  
266 diabetes. In addition, the T2DM cohort consisted mainly of individuals with no  
267 diabetic retinopathy (NDR), and a small proportion of individuals with diabetic  
268 retinopathy (DR) as defined by Early Treatment Diabetic Retinopathy Study (ETDRS)  
269 standards<sup>34</sup>.

### 270 **Image datasets and patient characteristics**

271 In the study, a large retinal fundus image dataset encompassing patient cohorts from  
272 the China Consortium of Fundus Image Investigation (CC-FII) was constructed,  
273 which consisted of cross-sectional datasets and longitudinal datasets. The

demographics and clinical information of the cohort participants are summarized in **Table 1** and **Extended Data Fig. 2**.

To develop an AI system for detecting CKD and T2DM, we first used a cross-sectional dataset (CC-FII-C) comprising 86,312 retinal fundus images from 43,156 participants, which are subsets from the CC-FII (**Fig. 1a** and **Table 1**). All subjects from the CC-FII-C were split randomly into mutually exclusive sets for training, tuning and “internal testing” of the AI algorithm at a 70%:10%:20% ratio. To evaluate the AI model’s generalization, we used two other independent retrospective population-based cohorts for external validation. The first external cohort consisted of non-selected 8,059 individuals who underwent an annual health-check from Guangdong Province (**Table 1**, see more details in **Methods**).

To further test its generalizability, we conducted a second external test (external test set 2), a population-based prospective study with retinal fundus images captured using a smartphone device. Given the potentially broad appeal of an AI-based medical diagnosis system based on non-invasive retinal imaging, we developed a low-cost hand-held smartphone camera attachment that would enable a healthcare professional or a patient to capture fundus images for assessment (**Extended Data Fig. 10**). The images could be uploaded to a Health Insurance Portability and Accountability Act (HIPAA)-compliant cloud service where the AI platform could autonomously grade incoming diagnostic requests. In this prospective study from September to November 2019, 3,081 patients were recruited in the COACS cohort with 6,162 smartphone-captured fundus images (see **Methods**).

To predict the development of CKD and T2DM, we used de-identified fundus images from two longitudinal datasets (**Extended Data Table 1**, see more details in **Methods**). The first dataset (CC-FII-L) for model development is a longitudinal cohort from CC-FII that contained 10,269 individuals from Tangshan City, Hebei Province, China, who underwent routine annual health-checks during a six-year follow-up period. The CC-FII-L dataset was randomly split into a developmental dataset and a longitudinal validation set (internal longitudinal test set) at a ratio of 8:2. For external validation, we used the second longitudinal dataset from Beijing, China (external longitudinal test set). This external longitudinal test set contained 3,376 individuals who had annual health check follow-ups for five years. The patient characteristics for each cohort can be found in **Extended Data Table 1**.

### Using the AI system to identify CKD and early CKD

We explored whether an AI algorithm could predict CKD’s presence or absence (including early or severe+ stages) in patients based on their fundus images and clinical metadata (for example, age, sex, height, weight, and blood pressure; see more details in **Methods**). Based on the definitions of CKD, we trained AI models to perform binary classification tasks. We first developed two models: baseline random

forest models using clinical metadata; and deep learning models using fundus images. We hypothesized that a model utilizing both clinical metadata and fundus images might be even more accurate, thus developing a combined AI model using both input modalities.

Training of an AI model using clinical metadata alone led to an area under the curve (AUC) of 0.861 on the receiver operating characteristic (ROC) curve (95% CI: 0.846-0.877) in the internal test set. In comparison, the AI model trained using fundus images alone produced a superior AUC of 0.918 (95% CI: 0.905-0.933). When trained using combined clinical metadata and fundus images, the combined AI model achieved comparable performance with an AUC of 0.930 (95% CI: 0.921-0.940) (**Fig. 2a**). We further validated these AI models using another independent external cohort (external test set 1) to demonstrate their generalizability. When tested on external test set 1, the AUC was 0.842 (95% CI: 0.827-0.856) for the clinical metadata-only model, 0.885 (95% CI: 0.873-0.899) for the fundus image model, and 0.898 (95% CI: 0.888-0.911) for the combined model (**Fig. 2b**).

Early CKD is defined by  $\text{eGFR} > 60 \text{ mL/min/1.73m}^2$  with albuminuria (a sign of kidney damage), whereas normal controls have  $\text{eGFR} > 60 \text{ mL/min/1.73m}^2$  without albuminuria. AI-based prediction of early CKD from normal controls (in the absence of information about albuminuria) is thus a useful means for early disease detection and prevention. The AI system's performance in predicting early CKD followed a similar trend across the three models (**Fig. 2d and 2e**). When evaluated using the internal test set from CC-FII-C, the clinical metadata model achieved an AUC of 0.805 (95% CI: 0.772-0.846), the fundus image model achieved 0.839 (95% CI: 0.805-0.868), and the combined model achieved 0.864 (95% CI: 0.837-0.894) (**Fig. 2d**). When tested on the external test set 1, the AUCs for predicting the presence of early CKD were 0.800 (95% CI: 0.780-0.824) for the clinical metadata model, 0.829 (95% CI: 0.811-0.849) for the fundus image model, and 0.848 (95% CI: 0.828-0.869) for the combined model (**Fig. 2e**).

Despite the AI model having been trained using images captured with standard hospital fundus cameras, the smartphone camera-captured images led to comparably good CKD detection performances based on fundus images alone. For the detection of CKD, the AI platform delivered an AUC of 0.817 (95% CI: 0.785-0.842) of the metadata-only model, an AUC of 0.870 (95% CI: 0.847-0.893) of the fundus-only model, and 0.897 (95% CI: 0.855-0.902) of the combined model in the external test set 2: "point-of-care" cohort (**Fig. 2c**). We further tested the AI performance for the staging of early CKD in the "point-of-care" study, which achieved an AUC of 0.787 (95% CI: 0.745-0.840), 0.834 (95% CI: 0.797-0.868), and 0.845 (95% CI: 0.812-0.892) for metadata-only, fundus-only, and combined model respectively (**Fig. 2f**). Together, these findings demonstrated not only the validity of the AI model, but also the potential real-life feasibility and utility of an AI-based retinal diagnosis



platform administered via personal mobile devices. This platform could assist medical professionals in screening and monitoring systemic microvascular diseases.

### **GFR prediction and CKD staging using retinal fundus images**

Estimated GFR (eGFR) is normally calculated based on measurement of serum creatinine level (see more details in Methods). We tested whether a patient's eGFR from measured serum creatinine could be predicted by an AI model using their fundus images alone. The level of agreement between the algorithm-predicted GFR and measured eGFR was assessed using a Bland-Altman plot. The AI model demonstrated good performance in eGFR prediction in the internal test set, achieving an intraclass correlation coefficient (ICC) of 0.65 (95% CI: 0.63-0.66) and mean absolute error (MAE) of 11.08 (95% CI: 10.80-11.35) (**Fig. 3a**). The model performed similarly well on the first external test set with an ICC of 0.62 (95% CI: 0.60-0.64) and an MAE of 12.91 (95% CI: 12.58-13.24) (**Fig. 3b**). When tested on the second external test cohort using smartphone captured images, the AI model achieved non-inferior performance with an ICC of 0.53 (95% CI: 0.50-0.55) and an MAE of 11.85 (95% CI: 11.42-12.27) (**Fig. 3c**). Bland-Altman plots showed a negative proportional bias (slope of linear fit); that is, the AI-based model underestimated the GFR level to a greater extent at high levels of GFR than at low levels. This proportional bias (slope of linear fit) was reduced after the model outputs were calibrated to have the same variance as that of the ground-truth measurements (Methods) (**Extended Data Fig. 3a-3c**).

Furthermore, a linear correlation between the AI-predicted GFR and measured eGFR showed strong associations with a coefficient of determination ( $R^2$ ) of 0.507 and Pearson's correlation coefficient (PCC) of 0.716 for the internal test set 1 (from CC-FII) (**Fig. 3d**). When evaluated on the external validation cohorts, the AI model achieved  $R^2$  of 0.481 and PCC of 0.700, and  $R^2$  of 0.327 and PCC of 0.577, for external test set 1 and 2, respectively (**Fig. 3e** and **3f**). These results suggest that the AI model was able to extract information predictive of eGFR, the key index of renal function, embedded subtly within fundus images.

Early detection and treatment of advanced to severe CKD could slow or prevent progression to end-stage renal failure and reduce mortality<sup>2</sup>. With the purpose of fundus imaging-based population screening of CKD in mind, we examined the AI model's performance in predicting the stage of CKD. A 'regression model' was first trained, generating a binary output in terms of presence or absence of severe+ CKD by applying a threshold after predicting the eGFR. Alternatively, a 'classification model' was also trained to perform the binary tasks of directly differentiating severe+ CKD from the other stages of CKD (early and advanced CKD), which achieved good performance with an AUC of 0.853 (95% CI: 0.799-0.891) in the internal validation dataset A (**Extend Data Fig. 3d**). Evaluated on the internal test set, the regression model showed comparable performance to the classification model in the detection of

403 severe+ CKD with an AUC of 0.825 (95% CI: 0.776-0.867) (**Extend Data Fig. 3d**).  
404 Similarly, when evaluated on the first external test set, the AUC for severe+ CKD  
405 identification was 0.842 (95% CI: 0.803-0.892) for the direct classification model and  
406 0.837 (95% CI: 0.788-0.877) for the indirect regression model (**Extend Data Fig. 3e**).

## 407 **Prediction of the development of CKD using longitudinal cohorts**

408  
409 Accurate prediction of CKD development has important implications for delivering  
410 targeted screening programs or risk-modifying intervention to those who would derive  
411 the most benefit from early disease detection. We hypothesized that not only could an  
412 AI algorithm predict the current CKD status of a patient based on their fundus images,  
413 but it may also be able to predict their future risk of CKD onset or progression. In this  
414 study, we implemented a deep learning AI model using baseline fundus-images and  
415 clinical metadata to predict risk of progression to CKD/advanced CKD in longitudinal  
416 cohorts (**Extended Data Table 1**).

417  
418 The performance of progression prediction models of CKD/advanced CKD using the  
419 Cox proportional hazards (CPH) model has been summarized in **Extended Data**  
420 **Table 4**. The metadata-based model achieved a C-index of 0.756 (0.699-0.810) on the  
421 internal test set and a C-index of 0.651 (0.569-0.730) on the external test set. When  
422 the deep learning features extracted from fundus images were combined with clinical  
423 metadata, the model performance improved to a C-index of 0.845 (0.789-0.910) on  
424 the internal test set and 0.719 (0.627-0.807) on the external test set. The combined  
425 progression prediction model had a statistically significant improvement in  
426 comparison to clinical metadata only-based prediction, as shown in **Extended Data**  
427 **Table 4** (Permutation test). In addition, we conducted analysis on the univariable and  
428 multivariable hazard ratio (HR) comparing the fundus-based models to known risk  
429 factors (**Extended Data Table 3**). We used HRs to examine and qualify the influence  
430 of specific factors on the rate of occurrence of a particular event rate (e.g., onset of  
431 disease) at a particular point in time. As shown in Extended Table 3, the adjusted HR  
432 of the fundus predictor is significant ( $p < 0.001$ ).

433  
434 Using the 6-year longitudinal data on CKD staging from our developmental patient  
435 cohorts (CC-FII-L), we further used the Kaplan-Meier method to stratify normal  
436 patients into three risk groups (the low, medium or high risks) for developing CKD or  
437 advanced to severe CKD (advanced+ CKD). The incidence of the CKD (per 1000  
438 person-years) stratified by risk groups of the AI model is shown in **Table 2**. For the  
439 Kaplan-Meier curves and log-rank tests, thresholds for the high-risk and the low-risk  
440 were based on the upper and lower quartiles of the predicted risk scores from the  
441 combined models in the developmental set. Significant separations of the low,  
442 medium and high-risk groups were also achieved in the internal test set ( $p < 0.001$ , **Fig.**  
443 **4a and 4c**). To assess the generalizability of the AI model, the same tests were  
444 performed on an external longitudinal test set of 3,215 patients with 5 years of

445 follow-up data. Once again, the AI model discriminated the low versus medium and  
446 high-risk groups for developing CKD or advanced+ CKD with high degrees of  
447 separation ( $p<0.001$  for both, **Fig. 4b and 4d**). The Kaplan-Meier curves for the risk  
448 stratification from the metadata-only model are illustrated in **Extended Data Fig.7**.  
449 **Extended Data Fig. 8** presents the cumulative hazards of three stratified risk groups  
450 in progressing to CKD/advanced CKD outcome at every time point. As shown, the  
451 combined model was discriminative in stratifying patients into low, medium and  
452 high-risk subgroups on both the internal longitudinal test set ( $p<0.001$ ) and the  
453 external longitudinal test set ( $p<0.001$ ).

454  
455 The prognostic accuracy of the AI system was assessed using time-dependent ROC  
456 analysis. The AUC for predicting the development of CKD was 0.844 (95% CI:  
457 0.787-0.888) in the internal test set. We further used the AUC at 4 years in the  
458 external longitudinal set to measure prognostic accuracy, which was 0.771 (95% CI:  
459 0.677-0.840) for predicting the onset of CKD (**Extended Data Fig. 9a and 9b**).

#### 460 **Using the AI system to identify T2DM**

461  
462 To further extend the scope of our AI system in the diagnosis and prediction of  
463 systemic microvascular diseases based on fundus images, we applied the same  
464 developmental framework to the detection of T2DM. The developmental dataset was  
465 divided into training, tuning and internal test sets (at a ratio of 7:1:2) to assess the  
466 models' performance (**Extended Data Table 1**).

467  
468 We first evaluated our models' performance in the detection of T2DM patients from  
469 normal controls. The AI system achieved an AUC of 0.828 (95% CI: 0.814-0.841) for  
470 the metadata-only model, an AUC of 0.923 (95% CI: 0.913-0.932) for the fundus  
471 image-only model, and an AUC of 0.929 (95% CI: 0.920-0.937) for the combined  
472 model in the internal test set (**Fig. 5a**). When evaluated on the first external test set,  
473 the model performed well, with an AUC of 0.796 (95% CI: 0.779-0.814) for the  
474 metadata-only model, 0.854 (95% CI: 0.839-0.871) for the fundus-only model, and  
475 0.871 (95% CI: 0.856-0.885) for the combined model (**Fig. 5b**). When evaluated on  
476 the second external test set using smartphone camera-captured images, the AI model  
477 achieved comparably good T2DM detection performance with an AUC of 0.762 (95%  
478 CI: 0.732-0.786) for the metadata-only model, 0.820 (95% CI: 0.788-0.853) for the  
479 fundus-only model, and 0.845 (95% CI: 0.822-0.869) for the combined model (**Fig.**  
480 **5c**).

481  
482 As a T2DM individual with DR can be readily detected using retinal images, we  
483 tested our AI performance independent of DR. Retinal fundus images of T2DM  
484 patients were divided into two subsets: (i) T2DM patients with diabetic retinopathy  
485 (DR) and (ii) T2DM patients with 'No Diabetic Retinopathy' (NDR) (**Extended Data**  
486 **Fig. 5a and 5b**). Comparable performance of the AI model in detecting T2DM with

DR or with NDR suggests that its accuracy is not heavily dependent on the presence of ETDRS-defined diabetic retinopathy. These results indicate that the AI model could detect T2DM based on fundus images prior to any apparent clinical manifestations of diabetic retinopathy (**Extended Data Fig. 5c**).

We further tested the ability of our models to predict the level of mean blood glucose from fundus images alone. Interestingly, the AI model achieved a relatively strong performance in the internal test set, the external test set 1 and the external test set 2, suggesting that blood glucose levels could be predicted and quantified from fundus images alone (**Extended Data Fig. 4a-4c**). Bland-Altman analysis and relating calibration were also performed to investigate the agreement between two approaches of the blood glucose measurement (**Extended Data Fig. 4d-4i**, see more details in **Methods**).

### **Prediction of the development of T2DM using longitudinal cohorts**

We next investigated the predictive performance of the AI system for the development of T2DM in normal patients over a 5-year period. The normal patients within the developmental cohort were stratified into the low, medium or high-risk groups for developing T2DM as defined by the upper and lower quartiles of our AI prediction. The incidence of the T2DM (per 1000 person-years) stratified by the three risk groups of the AI model is shown in **Table 2**. As shown in the Kaplan-Meier survival curve, a clean stratification of T2DM development rates between the low-risk and high-risk groups was observed in both the internal longitudinal test set ( $p < 0.001$ ) and the external longitudinal test set ( $p < 0.001$ ) (**Fig. 5c** and **5d**).

Subsequently, we performed progression analysis of the T2DM using CPH models. The performance of the metadata-based model and combined model is summarized in **Extended Data Table 5**. The metadata-based model achieved a C-index of 0.774 (0.732-0.819) on the internal test set and a C-index of 0.746 (0.706-0.775) on the external test set. When the deep learning features extracted from fundus images were combined with clinical metadata, the model performance improved to a C-index of 0.781(0.743-0.819) on the internal test set and 0.765 (0.723-0.799) on the external test set. We further performed univariable and multivariable survival analyses, including the basic prognostic factors, age, sex and height, in addition to the scores generated from the T2DM detection (**Extended Data Table 3**).

Among a total of 3,064 subjects with verified disease development outcomes, 185 (6.0%) developed T2DM within a 5-year follow-up period. The classical approach of ROC curve analysis considers an individual's event (disease) status as fixed over time; however, in practice, the disease status may change over time. In this study, we used time-dependent ROC curve analysis to assess the predictive ability of fundus image features. The AUC for predicting the onset of T2DM was 0.839 (95%CI: 0.799-0.886)

529 in the internal longitudinal test set and 0.824 (95%CI: 0.793-0.858) in the external  
530 longitudinal test set (**Extended Data Fig. 9c and 9d**).

### 531 **Visualization of evidence for CKD/T2DM prediction**

532  
533 Finally, to improve the interpretability of the AI model and shed light on its  
534 diagnostic mechanism, Integrated Gradients (IG) was used to generate saliency maps  
535 to highlight areas of the images that were important in determining the AI model's  
536 predictions. Several representative examples of original fundus images and their  
537 corresponding saliency maps are presented in **Fig. 6** and **Extended Data Fig. 5**.  
538 While microvascular pathologic changes are known to exist in CKD/T2DM, they may  
539 not be observable in fundus photographs, at least not to an ophthalmologist's eyes.  
540 The knowledge derived from saliency maps suggests that the AI model focuses on  
541 areas around the optic disc, central macula, retinal vessels, and other specific areas of  
542 abnormalities (e.g. in cases of diabetic retinopathy). These areas are commonly used  
543 by ophthalmologists to diagnose retinal diseases, suggesting that the AI model is  
544 learning clinically relevant features. Interestingly, the saliency maps for CKD and  
545 T2DM show somewhat distinct patterns of emphasis. In CKD, the predictive value  
546 locations appear to be the vessel branch points, arterial-venous junctions, and/or  
547 exudates or haemorrhages, suggesting that vascular health may play a role (**Fig. 6**). In  
548 T2DM, the highlighted 'points of interest' appear more scattered over the whole  
549 image, and in some instances, appear to correspond to the locations of DR features  
550 such as vascular tortuosity, venous dilatation, retinal haemorrhage and cotton wool  
551 spots (**Extended Data Fig. 5a and 5b**).

### 552 553 **Discussion**

554  
555 The common systemic microvascular diseases, CKD and T2DM, are major public  
556 health burdens and early detection is critical for successful clinical management and  
557 favorable outcomes. We hypothesized that the AI model might be able to diagnose  
558 CKD and T2DM by isolating and identifying subclinical changes that are  
559 undetectable by human observers. In this study, we developed an AI platform to  
560 detect CKD and T2DM with high degrees of accuracy and predict future disease  
561 development based on fundus images alone and prior to any clinical manifestation.  
562 These AI-derived diagnostic capabilities were not only applicable to clinical-grade  
563 retinal images obtained using professional cameras but also applicable to fundus  
564 images captured using smartphones. These findings could potentially provide a  
565 non-invasive, high throughput and low-cost screening tool for early detection of CKD  
566 and T2DM during the health screening. Additionally, the added value of fundus  
567 images in disease prognostication was further validated. The results raise the  
568 possibility of AI-based detection of other systemic diseases with retinal  
569 manifestations beyond clinicians' observational power.

570

571 Although the number, distribution and quality of the studies examining CKD  
572 prevalence and incidence have increased over the past decade, CKD surveillance  
573 capacity remains far less developed, with only 12.1% of patients identified as having  
574 CKD by primary care practitioners. Awareness of CKD remains low at 10% in US  
575 adults in part because CKD is usually a silent condition until its late stages<sup>35</sup>. In this  
576 regard, a recent report on CKD diagnosis using fundus images raises new possibilities  
577 for early disease detection<sup>21</sup>. However, earlier detection and prediction of disease  
578 development remains challenging but has great potential to help target interventions  
579 to individuals at high risk, thus reducing personal and economic burdens while  
580 improving the quality-of-life. Given the great prevalence of CKD, as well as an  
581 enormous screening demand, our AI system, which is capable of both diagnosing the  
582 disease and predicting the risk of onset, could improve CKD surveillance and capture  
583 of patients for early intervention. In this regard, the AI system should be viewed as a  
584 welcome addition to assist physicians and better direct healthcare resources.

585  
586 In addition, diabetes is a major global risk factor for developing high-risk CKD<sup>1</sup>.  
587 Good control of blood glucose levels reduces the risk of CKD and improves outcomes  
588 in patients with CKD<sup>36</sup>. Our retinal fundus image-based AI system showed its value  
589 for T2DM screening in general clinics and communities. As the majority of  
590 complications of diabetes are related to vascular damage, retinal fundus images  
591 provide highly relevant data for non-invasive and longitudinal assessment of vascular  
592 health. Though fasting blood glucose testing is relatively straightforward, the  
593 non-invasive fundus image-based screening system could broaden access and enable  
594 early detection and life-style modifications. In terms of monitoring for disease  
595 progression, compared with the HbA1c and random blood glucose testing,  
596 image-based analysis of vascular changes in the retina could provide alternative and  
597 objective means of assessing diabetic control and tissue damage.

598  
599 For real-world clinical applications, the operating point of an AI system could be set  
600 differently to balance the true positive rate (TPR) and the false-positive rate (FPR).  
601 Performance metrics for CKD/T2DM detection in the cohorts were determined by the  
602 operating points selected from the tuning dataset (shown in **Extended Data Table 6**  
603 **and 7**). We generated specific values of sensitivity, specificity, PPV and NPV of the  
604 AI system for systemic disease detection in order to meet various screening needs or  
605 clinical applications. To identify CKD cases with high confidence, we used a  
606 very-high decision threshold with a relatively higher PPV. For the external test set, we  
607 achieved a PPV of 88.4% (83.9-92.8) and 89.3% (80.8-95.5), while retaining a  
608 sensitivity of 34.8% (31.8-38.6) and 29.9% (21.5-34.3) (**Extended Data Table 6**). To  
609 identify normal controls with high confidence, we used a very-low decision threshold  
610 with a relatively higher NPV. For the external test set, we achieved a PPV of 99.7%  
611 (99.5-99.9) and 99.4% (98.5-100.0), while retaining a specificity of 37.5% (36.1-38.9)  
612 and 32.1% (29.0-35.5) (**Extended Data Table 6**).

613

In addition, this study has assessed the role of automated AI algorithms in the detection of CKD/T2DM using a low-cost smartphone-based imaging device. Camera-enabled smartphones are widely available around the world. The availability of retina specialists/trained retinal graders is, however, a major limitation in most countries. The portable and easy operation of our approach helps to provide screening in remote areas without the involvement of healthcare professionals. Therefore, the deployment of a smartphone-cloud based AI fundus diagnosis platform could provide access to essential diagnosis irrespective of regional resource variations and improve early detection of these treatable systemic diseases. In addition, this approach has potential as a non-invasive screening for multiple diseases in the general population and is not limited to patients with diabetes, as ocular imaging becomes more widely available. Such systems would also lead to significant cost-savings over traditional disease screening programs.

**This study has several differences to consider in comparison to previous studies<sup>20,21</sup>.**

Earlier literature in this area utilized a different definition of CKD. We used a more widely adopted definition of CKD based on the consensus clinical guideline: an eGFR  $\geq 60$  mL/min/1.73m<sup>2</sup> with albuminuria or eGFR  $< 60$  mL/min/1.73m<sup>2</sup>, while in a previous study for predicting CKD, the authors used a narrower clinical definition of CKD as a binary outcome on the basis of eGFR  $\leq 60$  mL/min per 1.73 m<sup>2</sup>. In contrast, in this study, CKD was defined by an eGFR  $\geq 60$  mL/min/1.73m<sup>2</sup> with albuminuria, or eGFR  $< 60$  mL/min, based on the clinical guidelines. Secondly, our time-to-critical event model based on longitudinal cohorts could provide a great utility in managing patients during their early disease course. Knowing which patients will progress to a CKD or advanced+ stages from a patient's first retinal fundus images will permit providers to triage and manage these patients while optimizing resources appropriately.

Moreover, as systemic diseases usually affect both eyes equally, we predicted with the AI system at the image-level, and averaged the image-level output at the patient-level for a systemic condition prediction. To confirm consistency between two eyes, we further conducted experiments comparing AI prediction performance based on the left eye only and right eye only data for CKD detection. This comparison showed a very strong inter-eye correlation (see more details in **Extended Data Fig. 3f**), supporting the accuracy and validity of the AI model.

Our study has several limitations which we hope to address in the future. First, since our AI was trained in a solely Chinese population, and tested in both an external Chinese cohort from several different geographic areas (the external test set 1&2), its generalizability in other racial populations will need to be further validated. Therefore, we added another external multi-ethnicity validation cohort (**Extended Data Table 8**) consisting of individuals from Kashi (also called Kashgar) in Xinjiang Autonomous Region (Uighur ethnicity) and Macau (Portuguese ethnicity), in which the AI model showed good performance (**Extended Data Fig. 6**). Additional training on more

diverse clinical and demographic cohorts may further improve diagnostic accuracy and clinical utility in a broad range of populations. Another limitation is related to the measurement of eGFR; inaccurate estimation may exist in patients with rapid deterioration of renal function. In this study, patients with combined AKI diagnosis (or a low incidence of AKI) were excluded to minimize this interference. Thus, such a calculation strategy can still reflect the characterization of the patient's renal function progression. In addition, whether additional clinical metadata (such as blood pressure trends, smoking status, alcohol consumption level, and family history) could further improve the accuracy of the predictions need to be explored. While the model appears well suited for diagnostic screening, it remains limited in its ability to provide prognostic information to individual patients or insights into pathogenic mechanisms based on saliency maps.

## Methods

### Dataset characteristics

To develop an AI system for the detection of CKD and T2DM, fundus image data were collected from the China Consortium of Fundus Image Investigation (CC-FII), which included the following participants: China suboptimal health cohort study (COACS) in Tangshan City, Hebei Province; the Peking University First-Affiliated Hospital and the Peking University Second-Affiliated Hospital, both in Beijing; West China Hospital in Chengdu, Chongqing Renji Hospital in Chongqing, both in Sichuan Province; Kunshan Hospital of Jiangsu University, Kunshan, Jiangsu Province. All procedures were performed as a part of a patients' routine annual health check. Institutional Review Board (IRB)/Ethics Committee approvals were obtained in all locations and all participating subjects signed a consent form.

The COACS is a community-based, prospective study to investigate how suboptimal health status contributes to the incidence of NCD (non-communicable chronic diseases) in Chinese adults<sup>37</sup>. This COACS study has two phases, a cross-sectional survey followed by a longitudinal study. The participants were recruited from Tangshan city, which is a large, modern industrial city and adjoins two mega cities: Beijing and Tianjin. In phase I, all participants underwent clinical laboratory measurements. In the second phase, a long-term yearly clinical follow-up has been performed until 2024, with the purpose of better understanding how suboptimal health, environmental and genetic risk factors contribute to the development of major chronic diseases. We have elected to use this cohort for our study because it balances healthy subjects and those with major chronic diseases such as CKD and T2DM.

For the CKD and T2DM detection, we used retinal fundus images from retrospective datasets. The first one, a cross-sectional cohort (CC-FII-C) from CC-FII, included a total of 86,312 fundus images from 43,156 subjects as the developmental dataset of



our AI models for systemic disease detection. All subjects from the developmental cohort were split into mutually exclusive sets for a training set and a tuning set (70%:10%), as well as an internal test set (20%). Detailed participant characteristics are shown in **Table 1**. The external test set 1 included subjects undergoing an annual health check in Zhong Shan Ophthalmic Center of Sun Yat-sen University and Nanfang Hospital both in Guangzhou, Guangdong Province. Study participants also subsequently underwent ophthalmological examinations with fundus imaging. Only retinal fundus images were included in this study with 8,059 participants from the external test set 1. We further tested our AI system in the external test set 2, which is a prospective “point-of-care” settings, for the evaluation of the generalizability of the AI system with smartphone-based device. For this prospective study, we enrolled a total of 3,081 subjects, including 228 CKD patients, and 1,189 T2DM patients in Tangshan from September 16, 2019, to November 15, 2019. These participants were a part of the COACS study, yet they were independent of the development dataset. To demonstrate the generalizability of the AI model in other ethnic groups, we added another external multi-ethnicity validation cohort consisting of 400 individuals from the First Regional Hospital of Kashi in Xinjiang Autonomous Region; as well as the University Hospital of Macau University of Science and Technology and Hospital Kiang Wood, both in Macau, China.

To develop an AI system for the prediction of the incidence of CKD and T2DM, we further used two longitudinal cohorts (**Extended Table 1**). The first longitudinal dataset CC-FII-L is a subset from CC-FII, which consisted of 10,269 participants for routine annual health checks with a follow-up period of 6 years. All the subjects from the CC-FII-L were randomly divided for developmental set and internal longitudinal test set (3,376 individuals), in an 8:2 ratio. Another longitudinal dataset, an external longitudinal test set, was used for external validation of the AI system for incidence prediction of the systemic diseases. This is a population-based study of Chinese from Beijing, China, which recruited patients from hospitals or health centers for annual health checks, including diabetic retinopathy screening in Peking University’s First Affiliated Hospital and Third Affiliated Hospital. A total of 3,376 subjects were used as external validation (external longitudinal test set) of our AI models for the incidence prediction of CKD and T2DM. The patient characteristics for each cohort can be found in **Extended Data Table 1**.

#### **Image quality control**

The retinal fundus images were captured using a variety of standard fundus cameras, including Topcon TRC-NW6 (Topcon Corp, Tokyo, Japan), Zeiss Visucam 224 (Carl Zeiss Meditec AG, Jena, Germany), Canon CR6-45NM (Canon Inc, Tokyo, Japan), and KOWA Nonmyd  $\alpha$ -DIII (Kowa Company Ltd, Aichi, Japan). During the image grading process, all fundus images were first de-identified to remove any patient-related information. Participants were recruited from an annual health check population that underwent a physical examination by a physician and fundus image

photography. For screening and grading retinal fundus images of diabetic retinopathy, a hierarchical two-tier grading process was performed by ten phase I and five phase II graders. Phase I graders consisted of individuals trained by ophthalmologists and evaluated to perform at least 95% accuracy determined by a 1000-fundus-image quiz of various retinal diseases. Phase II graders consisted of ophthalmologists who individually reviewed every image classified by phase I graders. To check consistency among phase II graders, 20% of images were randomly selected and reviewed by 3 senior retinal specialists. The second tier of five ophthalmologists independently read and verified the true labels for each image. In order to account for disagreement, the evaluation test set was also checked by expert consensus. About 9% of the study participants were excluded due to poor photographic quality/unreadable images and missing clinical diagnosis. After establishing the consensus diagnoses, images were transferred to the AI team to develop a deep learning algorithm for image-based classification.

#### **Definition and criteria for disease diagnosis**

The following criteria were used to define the systemic disease and related disease staging. According to the international kidney disease clinical classification based on the KDIGO guideline (Kidney Disease: Improving Global Outcomes (KDIGO) Clinical Practice Guideline), the staging and actual risk of adverse outcomes of kidney disease are stratified by the renal glomerular filtration function defined as the glomerular filtration rate (GFR) categories<sup>38</sup>. GFR is an indicator of overall kidney function, which equals the total amount of fluid filtered through the functioning nephrons per unit of time<sup>39</sup>. The estimated GFR (eGFR) in this study was based on the Chronic Kidney Disease Epidemiology Collaboration (CKD-EPI) Equation. This equation has been extensively validated in Chinese and Asian populations<sup>40,41</sup>.

In clinical practice, CKD was diagnosed as an eGFR of more than 60 mL/min/1.73 m<sup>2</sup> with albuminuria or less than 60 mL/min/1.73 m<sup>2</sup>, confirmed in at least two visits separated by three months. Once a diagnosis of CKD has been made, the next step is to determine staging, which is determined by the extent to which GFR has decreased. Staging of kidney function is classified with eGFR as follows: Stage 1 (eGFR > 90 mL/min/1.73 m<sup>2</sup>, with albuminuria), Stage 2 (eGFR 60-89 mL/min/1.73 m<sup>2</sup> with albuminuria), Stage 3 (eGFR 30-59 mL/min/1.73 m<sup>2</sup>), and Stage 4+ (eGFR < 30 mL/min/1.73 m<sup>2</sup>). For the CKD stage classification, we defined the potential or mild kidney impairment, called early CKD, corresponding to Stages 1 and 2. Stage 3 is denoted as advanced CKD. In addition, the detection and early intervention of CKD severe+ are also crucial, which is defined with an eGFR cutoff with 30 mL/min/1.73 m<sup>2</sup> (corresponding to stage 4 or above). Normal controls were defined as eGFR above 60 mL/min/1.73 m<sup>2</sup> without albuminuria, checked by a negative urine dip-stick test. In our study, we utilized the images and corresponding eGFR measurements of already-diagnosed CKD patients. For a patient with multiple visits, we use the labels

784 of retinal fundus images taken corresponding to the visit when first establishing the  
785 diagnosis.

786

787 T2DM was diagnosed by fasting blood glucose  $\geq 7.0$  mmol/L at least 2 times, or as an  
788 HbA1c value of 6.5 % or more, and/or by a history of drug treatment for diabetes. For  
789 the detection of T2DM patients, there exist two subsets of fundus images: a diabetic  
790 retinopathy (DR) group and a ‘diabetes mellitus with no DR’ group (NDR) group. In  
791 the study, we conducted an experiment of NDR versus normal controls to investigate  
792 the performance of detecting asymptomatic T2DM patients. We also compared the  
793 performance with/without DR, and the results show that the AI system could learn  
794 early eye features of T2DM patients and showed comparable results with/without DR.  
795 In our detection study, we utilized the images and corresponding clinical data,  
796 including laboratory values of already-diagnosed T2DM patients.

## 797 **Deep learning and transfer learning methods**

### 798 *Categories of tasks for systemic diseases detection*

799 To develop a diagnostic platform which is capable of detecting vascular-related  
800 systemic diseases based on retinal images, we trained separate deep learning models  
801 for each task. More specifically, we had two types of prediction tasks, including a  
802 "regression" task and a "classification" task. For the regression tasks, we trained two  
803 models to predict continuous values of eGFR and fasting blood glucose, respectively.  
804 The detection of systemic disease or staging was treated as a classification task. We  
805 also performed the binary classification for T2DM detection. For each of these tasks,  
806 including CKD detection and T2DM detection, we compared the performance of three  
807 prediction models, each with a different set of input data. As a baseline, we used  
808 metadata-only models., which comprised age, sex, blood pressure, hypertension,  
809 height, weight and BMI, to develop random forest and logistic regression classifiers.  
810 We included diabetes as a covariate when building our AI model in CKD prediction.  
811 We further used fundus-only models based on deep convolutional neural networks  
812 (CNNs) with fundus images. Finally, we trained combined models which integrate  
813 fundus image data and clinical information. The image feature vector derived from  
814 the CNN model was concatenated with the clinical feature of the same patient. A  
815 multi-layer perception (MLP) took these features as input for classification.

816 For training and tuning the random forest model, the number of trees in the forest  
817 "n\_estimators" was with the default setting, and the maximum number of leaf nodes  
818 "max\_leaf\_nodes" and the minimum number of samples required to be at a leaf node  
819 "min\_samples\_leaf" were tuned by parameter search. Finally, they were set with  
820 n\_estimators=100, max\_leaf\_nodes=31 and min\_samples\_leaf=20.” Similarly, we  
821 added logistic regression models for the CKD and T2DM detection, with regularized  
822 likelihood estimation (L2 penalty) and regularization coefficient C=1.0 (Extended  
823 Data Table 2).

## 824 *Fundus image enhancement*

825 To capture the non-specific anatomical and physiological features on fundus images  
826 relevant to systemic diseases, we proposed image enhancement to improve the  
827 performance of the AI models. Two methods were utilized for fundus image  
828 enhancement, including Contrast Limited Adaptive Histogram Equalization  
829 (CLAHE)<sup>42</sup> and color normalization<sup>43</sup>. CLAHE enhancement is conducted by dividing  
830 the image into local regions and applying histogram equalization over all  
831 neighborhood pixels. Specifically, we converted the input fundus images from RGB  
832 color space into LAB color space. After applying the CLAHE on the lightness  
833 channel, we then converted it back to RGB. Compared with the original fundus  
834 images, CLAHE algorithm enhanced the details and visibility level of the fundus  
835 image. The fundus image normalization method was performed as follows:  
836  $x' = \alpha x - \alpha \text{Gauss}(x, \mu, \Sigma, s \times s) + \beta$ , where  $x$  is the input image,  $x'$  is the normalized  
837 image,  $\alpha$  and  $\beta$  are parameters, and  $\text{Gauss}(x, \mu, \Sigma, s \times s)$  is the Gaussian filtered  
838 image with a Gaussian kernel  $(\mu, \Sigma)$  of size  $s \times s$ . We used  $\alpha = 4$  and  $\beta = 128$ ,  
839  $\Sigma = I$ , and  $s = 10$ , following the settings of Liu et al<sup>43</sup>. With image normalization, we  
840 could reduce the color scope bias among fundus images taken under different lighting  
841 or device conditions. The effectiveness of the image enhancement method can be seen  
842 in **Extended Data Fig. 1**.

## 843 *Model development and training*

844 Convolutional neural networks (CNNs) were used to analyze and classify the fundus  
845 images in this study. With the ResNet-50<sup>44</sup> as the backbone, we pretrained on the  
846 ImageNet dataset for all deep learning models demonstrated. ResNet-50 is a  
847 five-stage network with one convolution and four identity blocks, which utilizes skip  
848 connections to overcome the degradation problem of deep learning models. For  
849 "regression" tasks of continuous values prediction (eGFR and fast blood glucose), a  
850 fully connected layer with one scalar as output was used as the final layer in the  
851 ResNet-50 model. For binary classification tasks, an additional softmax layer besides  
852 a fully connected layer was attached to the model. Retraining consisted of loading the  
853 convolutional layers with pretrained weights, newly initializing additional layers for  
854 our regression and binary classification tasks and training models on the  
855 corresponding development sets.

856 We used a three-layer MLP for combined models. Each of the two hidden layers has  
857 128 nodes and was applied with the ReLU activation function. The MLP was jointly  
858 trained with the CNN. The Mean-Square Error (MSE) loss was used as an objective  
859 function for the "regression" tasks of prediction of continuous values and the Binary  
860 Cross Entropy (BCE) loss was used for binary "classification" tasks. Training of  
861 models by back-propagation of errors was performed in batches of 32 images resized  
862 to  $512 \times 512$  pixels for 50 epochs with a learning rate of  $10^{-3}$ . Training was  
863 performed using the Adam optimizer with a weight decay of  $10^{-6}$ . Transformations  
864 of random horizontal and vertical flip were added to each batch during training as data

865 augmentation in order to enable an improved and generalized network learning. The  
866 models were implemented using PyTorch<sup>45</sup>. We randomly divided the developmental  
867 dataset into a training set (7/8 of the development set) and a tuning set (1/8 of the  
868 development set) to develop our model. The models selected for evaluation on  
869 validation sets were the models with the best validation loss on the tuning set. There  
870 were no samples overlapping at the patient level in training and evaluation sets.

### 871 ***Model Ensemble***

872 To improve the overall performance of the AI, we applied a model ensemble. For  
873 each task, we trained four model instances with different processed fundus images as  
874 input. Each input image was pre-processed into three variations by applying CLAHE  
875 only, normalization only, and both CLAHE and normalization. Thus, for each task,  
876 we had four models with the same architecture trained in parallel on the same  
877 development set but with each using differently pre-processed fundus images. The  
878 reported predictions were obtained by averaging the outputs of the four model  
879 instances.

### 880 **Prediction of the development of systemic diseases using longitudinal cohorts**

881  
882 For the incidence analysis of the CKD, we denoted the index data as the time when  
883 the subjects were without CKD at baseline. All participants who underwent a urine  
884 analysis with a negative result at the baseline visit using urine dip-stick test, were  
885 included for the CKD incidence analysis. The incidence data were denoted as the time  
886 when the subjects were recorded as having CKD/ advanced+ CKD during the  
887 follow-up visits. Similarly, we predicted the development of T2DM, with the index  
888 data defined as one without T2DM at the first visit. The development of T2DM was  
889 diagnosed as T2DM incidence data (or end-point) within the yearly clinical follow-up.  
890 We trained the Cox' proportional hazard (CPH) models on the training and tuning  
891 set  
892 using variables based on the metadata and fundus image, which comprised age, sex,  
893 blood pressure, height, weight, BMI, hypertension, T2DM, and predicted z-score  
894 (standard score) of the first visit generated from the CKD/T2DM detection model.  
895

896 According to the risk scores of the first visit from the CPH model for the CKD/T2DM  
897 detection, the patients are triaged into three groups: low, medium, and high risk  
898 according to the upper and lower quartiles of predicted risk scores in the tuning set,  
899 respectively. Kaplan-Meier curves were constructed for the risk groups, and the  
900 significance of differences between group curves was computed using the log-rank  
901 test. Time-dependent ROC curves<sup>46</sup> were used to quantify model performance on  
902 validation sets at the time of interest. ROC curves were constructed at a landmark  
903 time from predicted risk scores of relative patients made using the model. The  
904 univariable and multivariable CPH models were fitted. Two multivariable CPH  
905 models were developed, a combined metadata and fundus model and a metadata-only

906 model serving as a baseline model. Statistical significance of hazard ratios and  
907 adjusted hazard ratios of CPH models were evaluated using the likelihood ratio test.

### 908 **Handheld fundoscopy using a smartphone attachment**

909  
910 A 22-diopter double-convex aspheric condensing lens (Volk Optics, Ohio, USA) was  
911 mounted inside a custom-designed smartphone fundoscope attachment (**Extended**  
912 **Data Fig. 10**). The plastic components of the fundoscope were computer-designed on  
913 SolidWorks with collision simulation, converted to Standard Triangle Language  
914 format, and 3D printed (fused filament fabrication) in polylactic acid to a resolution of  
915 100 microns. The lens was stably anchored to the fundoscope cone using a rubber ring  
916 and printed lens locking system. Two perpendicular polarizing filters were also  
917 incorporated within the optical system to minimize the reflection of the smartphone  
918 flashlight from the cornea.

919  
920 We calculated the cone-shaped offset of the iPhone (Apple Inc, USA) from the  
921 condensing lens and the distance between the condensing lens and the anterior  
922 principal plane of the eye using 12.3 cm as the focal length of the condensing lens.  
923 Thus, the condensing lens worked in harmony with the iPhone camera light source  
924 and the optical system of the eye to achieve Maxwellian illumination of the retina and  
925 project an in-focus and widefield image onto the camera sensor. Fundus images were  
926 captured using this smartphone-mounted fundoscope in a prospective study within the  
927 COACS study. Informed consent was obtained from patients prior to pupil dilation  
928 and retinal photography using the standard operating procedure below. The same AI  
929 system for detecting CKD or T2DM as used for analyzing professional fundus  
930 camera-derived images was used for the handheld fundoscopy-derived images, which  
931 detected systemic diseases of CKD or T2DM. The performance of the model was  
932 evaluated using ROC curves.

### 933 ***Imaging protocol using the smartphone attachment***

934  
935 Standard operating procedure for fundus image capture using the smartphone  
936 fundoscope attachment:

- 937 1. The pupil is dilated with a drop of 1% tropicamide.
- 938 2. Select 'New Patient' from the main program display and enter patient information.
- 939 3. First, hold the iPhone X (Apple Inc, Cupertino, CA, USA) with fundoscope  
940 attachment in the right hand approximately 10 cm from the patient's eye to obtain  
941 a red reflex at the center of the display.
- 942 4. Then, slowly move the imaging device towards the eye while keeping the red  
943 reflex centered on the display. When the red reflex fills the entire field-of-view,  
944 stabilize the end of the fundoscope attachment with the left hand, which leans on  
945 the patient's forehead for stability.



5. Make fine adjustments to the distance between the end of the fundoscope attachment and the eye to obtain a focused image of the retina. Press the capture button on the iPhone to acquire the image.
6. Base on the quality of images obtained, the photographer may acquire up to 9 images per eye, with the best quality image being selected for the study.

## Interpretation of AI predictions

The difficulty in obtaining clinically intelligible features remains the greatest drawback of artificial neural network-based systems. A visualization tool is needed that would enable clinicians to understand important clinical variables in real time. To this end, we employed the Integrated Gradient algorithm<sup>47</sup> in order to produce “visual explanations”. Gradient-based visualization methods use gradients to quantify the importance of each pixel in the image. The importance of each pixel in the image to the correct predictions of the models can be quantified by  $\frac{\partial y_i}{\partial x}$ , where  $y_i$  is the model output for class  $i$  and  $x$  is the input image. However, gradient saturation presents a problem where the gradients of the output with respect to the input can be small even though they are important for the model output. Such phenomena can happen when the model outputs for the correct class reach a certain magnitude. To overcome this problem, the Integrated Gradient method improve the measurement of importance as follows:  $I_i(x) = (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial f(x' + \alpha(x - x'))}{\partial x_i} d\alpha$ , where  $I_i$  is the integrated gradients for pixel  $i$ ,  $x$  is the input image,  $x'$  is the baseline image,  $x_i$  and  $x'_i$  are values of pixel  $i$  in  $x$  and  $x'$ ,  $f$  is the model to be visualized. The saliency maps generated by integrated gradient indicate the effect of each pixel on the model predictions. We applied Gaussian filtering to saliency maps for smoothness on the fundus images.

## Statistical analysis

To evaluate the performance of regression models for continuous values prediction in this study, we calculated Mean Absolute Error (MAE), R-square (R<sup>2</sup>), and Pearson Correlation Coefficient (PCC). We applied the Bland-Altman plot<sup>48</sup> to display the difference between the measured eGFR and the predicted value of a sample against the average of the two. With 95% limits of agreement and Intraclass Correlation Coefficient (ICC), we evaluated the agreement of the predicted eGFR and actual eGFR. The models' performance on binary classification predictions was evaluated by Receiver Operating Characteristic (ROC) curves of sensitivity versus 1 – specificity. The Area Under the Curve (AUC) of ROC curves were reported with 95% Confidence Intervals (CIs). The 95% CIs of AUC were estimated with the non-parametric bootstrap method (1,000 random resampling with replacement). Sensitivity and specificity were determined by the selected thresholds on the validation set. The prediction of continuous values of eGFR and blood glucose level were evaluated with regression models. The CKD and T2DM detection were evaluated with binary classification models. For each subject, we made a prediction with the AI system at the image-level, and then averaged the image-level output at a

987 patient-level for a final prediction in each patient. ICC was used to assess the  
988 agreement between AI predicted values of left and right eyes, where stage CKD was  
989 measured by predicted eGFRs and diabetes, measured by log-likelihood ratios of  
990 predicted probability of T2DM presence. We calculated incidence rate for the whole  
991 cohort and for each risk group by three-strata shown as the number of events per  
992 1,000 person-years at risk. The Byar Poisson approximation method was used to  
993 calculate 95% CIs of incidence<sup>49</sup>. Then Kaplan-Meier estimators were constructed for  
994 different risk groups, and the significance of differences between groups was tested  
995 by log-rank tests. We used the time-dependent AUC (area under the curve) at four  
996 years and five years to measure model performance. The Kaplan-Meier curve and the  
997 time-dependent ROC-AUC were calculated using the Python packages of lifelines  
998 (version 0.25.5) and scikit-survival (version 0.14.0).

#### 999 **Data availability statement**

1000 To protect patient confidentiality, we have deposited patient retinal image data in a  
1001 secured and patient confidentiality compliant cloud in China in concordance with data  
1002 security regulations. Data access can be requested by writing to the corresponding  
1003 authors. All data and code access requests will be reviewed and (if successful) granted  
1004 by the Data Access Committee.

#### 1005 **Code availability statement**

1006  
1007 The deep learning models were developed and deployed using standard model  
1008 libraries and the PyTorch framework. Custom codes were specific to our development  
1009 environment and used primarily for data input/output and parallelization across  
1010 computers and graphics processors. We will make our codes available upon request  
1011 and approval by a Data Access Committee.

#### 1012 **Acknowledgements**

1013  
1014 This study was funded by the National Key Research and Development Program of  
1015 China (2019YFB1404804, 2017YFC1104600, 2017YFC0112402), Guangzhou  
1016 Regenerative Medicine and Health Guangdong Laboratory, National Natural Science  
1017 Foundation of China (61906105, 61872218 and 61721003), Macau University of  
1018 Science and Technology, Tsinghua University Initiative Scientific Research Program.  
1019 We thank members of Zhang, Yuan, and Wang groups for their assistance. We thank  
1020 many volunteers and physicians for grading retinal photographs.

#### 1021 **Author Contributions**

1022 J.X., J.Y., ZH.L., WJ. C., WQ. Xu., X. L., A. O., GX. Z., LH. Z., C. Z., OL. L., E. Z.,  
1023 J. Z., SQ. H., KM. X., W. W., S. N., M. L., J. Z., M. P. V., M. A., JQ. W., A. W., XG.  
1024 Z., Q. Z., P.T., KX.Z., T.X., GY. W., and K. Z. collected and analyzed the data. K.Z.  
1025 and GY.W. conceived and supervised the project. J.Y-N.L and D. L. help with the



1026 data analysis and interpretation. K.Z., KM.X., GY.W. wrote the manuscript. All  
1027 authors discussed the results and reviewed the manuscript.

1028

## 1029 **Competing Interests statement**

1030

1031 The authors declare no competing financial interests.

1032

## 1033 **References**

1034

- 1035 1. Collaboration, G.B.D.C.K.D. Global, regional, and national burden of chronic  
1036 kidney disease, 1990-2017: a systematic analysis for the Global Burden of  
1037 Disease Study 2017. *Lancet* **395**, 709-733 (2020).
- 1038 2. Levin, A., *et al.* Global kidney health 2017 and beyond: a roadmap for closing  
1039 gaps in care, research, and policy. *Lancet* **390**, 1888-1917 (2017).
- 1040 3. Kooman, J.P., Kotanko, P., Schols, A.M., Shiels, P.G. & Stenvinkel, P.  
1041 Chronic kidney disease and premature ageing. *Nature reviews. Nephrology* **10**,  
1042 732-742 (2014).
- 1043 4. Saeedi, P., *et al.* Global and regional diabetes prevalence estimates for 2019  
1044 and projections for 2030 and 2045: Results from the International Diabetes  
1045 Federation Diabetes Atlas, 9(th) edition. *Diabetes research and clinical  
1046 practice* **157**, 107843 (2019).
- 1047 5. Leasher, J.L., *et al.* Global Estimates on the Number of People Blind or  
1048 Visually Impaired by Diabetic Retinopathy: A Meta-analysis From 1990 to  
1049 2010. *Diabetes care* **39**, 1643-1649 (2016).
- 1050 6. Balakumar, P., Maung, U.K. & Jagadeesh, G. Prevalence and prevention of  
1051 cardiovascular disease and diabetes mellitus. *Pharmacological research* **113**,  
1052 600-609 (2016).
- 1053 7. From the Center of Disease Control and Prevention. Lower extremity  
1054 amputation episodes among persons with diabetes--New Mexico, 2000. *Jama*  
1055 **289**, 1502-1503 (2003).
- 1056 8. American Diabetes, A. 11. Microvascular Complications and Foot Care:  
1057 Standards of Medical Care in Diabetes-2020. *Diabetes care* **43**, S135-S151  
1058 (2020).
- 1059 9. Luk, A.O., *et al.* Quality of care in patients with diabetic kidney disease in  
1060 Asia: The Joint Asia Diabetes Evaluation (JADE) Registry. *Diabetic medicine :  
1061 a journal of the British Diabetic Association* **33**, 1230-1239 (2016).
- 1062 10. Wu, B., Zhang, S., Lin, H. & Mou, S. Prevention of renal failure in Chinese  
1063 patients with newly diagnosed type 2 diabetes: A cost-effectiveness analysis.  
1064 *Journal of diabetes investigation* **9**, 152-161 (2018).
- 1065 11. Esteva, A., *et al.* A guide to deep learning in healthcare. *Nature medicine* **25**,  
1066 24-29 (2019).
- 1067 12. Norgeot, B., Glicksberg, B.S. & Butte, A.J. A call for deep-learning healthcare.  
1068 *Nature medicine* **25**, 14-15 (2019).
- 1069 13. Ravizza, S., *et al.* Predicting the early risk of chronic kidney disease in  
1070 patients with diabetes using real-world data. *Nature medicine* **25**, 57-59  
1071 (2019).
- 1072 14. Topol, E.J. High-performance medicine: the convergence of human and  
1073 artificial intelligence. *Nature medicine* **25**, 44-56 (2019).

- 1074 15. De Fauw, J., *et al.* Clinically applicable deep learning for diagnosis and  
1075 referral in retinal disease. *Nature medicine* **24**, 1342-1350 (2018).
- 1076 16. Kermany, D.S., *et al.* Identifying Medical Diagnoses and Treatable Diseases  
1077 by Image-Based Deep Learning. *Cell* **172**, 1122-1131 e1129 (2018).
- 1078 17. Liang, H., *et al.* Evaluation and accurate diagnoses of pediatric diseases using  
1079 artificial intelligence. *Nature medicine* **25**, 433-438 (2019).
- 1080 18. Wang, C., Elazab, A., Wu, J. & Hu, Q. Lung nodule classification using deep  
1081 feature fusion in chest radiography. *Computerized Medical Imaging and*  
1082 *Graphics* **57**, 10-18(2017).
- 1083 19. Poplin, R., *et al.* Prediction of cardiovascular risk factors from retinal fundus  
1084 photographs via deep learning. *Nature biomedical engineering* **2**, 158-164  
1085 (2018).
- 1086 20. Rim, T.H., *et al.* Prediction of systemic biomarkers from retinal photographs:  
1087 development and validation of deep-learning algorithms. *The Lancet Digital*  
1088 *Health* **2**, e526-e536 (2020).
- 1089 21. Sabanayagam, C., *et al.* A deep learning algorithm to detect chronic kidney  
1090 disease from retinal photographs in community-based populations. *The Lancet*  
1091 *Digital Health* **2**(2020).
- 1092 22. Liu, T.Y.A. Smartphone-Based, Artificial Intelligence-Enabled Diabetic  
1093 Retinopathy Screening. *JAMA Ophthalmology* **137**, 1188-1189 (2019).
- 1094 23. Chen, C.-F., Lee, G.G., Sritapan, V. & Lin, C.-Y. in 2016 IEEE International  
1095 Workshop on Signal Processing Systems (SiPS). 130-135 (*IEEE*) (2016).
- 1096 24. Howard, A.G., *et al.* MobileNets: Efficient Convolutional Neural Networks for  
1097 Mobile Vision Applications. *arXiv preprint arXiv* **1704.04861**(2017).
- 1098 25. Wu, Y., Lim, J. & Yang, M.H. Object Tracking Benchmark. *IEEE*  
1099 *transactions on pattern analysis and machine intelligence* **37**, 1834-1848  
1100 (2015).
- 1101 26. Schroff, F., Kalenichenko, D. & Philbin, J. in Proceedings of the IEEE  
1102 conference on computer vision and pattern recognition. **815-823**.
- 1103 27. Bahdanau, D., Cho, K. & Bengio, Y. Neural Machine Translation by Jointly  
1104 Learning to Align and Translate. *ArXiv* **1409**(2014).
- 1105 28. Mikolov, T., Chen, K., Corrado, G.s. & Dean, J. Efficient Estimation of Word  
1106 Representations in Vector Space. *Proceedings of Workshop at ICLR*  
1107 **2013**(2013).
- 1108 29. Ignatov, A., *et al.* *AI Benchmark: Running Deep Neural Networks on Android*  
1109 *Smartphones*, (2018).
- 1110 30. Organization, W.H. Density of physicians (total number per 1000 population,  
1111 latest available year). *Last accessed Oct* [http://www. who.](http://www.who.int/gho/health_workforce/physicians_density/en)  
1112 [int/gho/health\\_workforce/physicians\\_density/en](http://www.who.int/gho/health_workforce/physicians_density/en). (2016).
- 1113 31. Poushter, J. Smartphone ownership and internet usage continues to climb in  
1114 emerging economies. *Pew Research Center* **22**, 1-44 (2016).
- 1115 32. Vos, T., *et al.* Global, regional, and national incidence, prevalence, and years  
1116 lived with disability for 310 diseases and injuries, 1990 - 2015: a systematic  
1117 analysis for the Global Burden of Disease Study 2015. *The Lancet* **388**,  
1118 1545-1602 (2016).
- 1119 33. Gansevoort, R.T., *et al.* Lower estimated GFR and higher albuminuria are  
1120 associated with adverse kidney outcomes in both general and high-risk  
1121 populations. *Kidney International Supplements* **80(1):93-104**(2011).

1122 34. Group, E.T.D.R.S.R. Grading diabetic retinopathy from stereoscopic color  
1123 fundus photographs—an extension of the modified Airlie House classification:  
1124 ETDRS report number 10. *Ophthalmology* **98**, 786-806(1991).

1125 35. Tuot, D.S., *et al.* Chronic kidney disease awareness among individuals with  
1126 clinical markers of kidney dysfunction. *Clinical journal of the American*  
1127 *Society of Nephrology : CJASN* **6**, 1838-1844 (2011).

1128 36. Tuttle, K.R., *et al.* Diabetic kidney disease: a report from an ADA Consensus  
1129 Conference. *American journal of kidney diseases : the official journal of the*  
1130 *National Kidney Foundation* **64**, 510-533 (2014).

1131 37. Wang, Y., *et al.* China suboptimal health cohort study: rationale, design and  
1132 baseline characteristics. *Journal of translational medicine* **14**, 291 (2016).

1133 38. Levin, A., *et al.* Kidney disease: Improving global outcomes (KDIGO) CKD  
1134 work group. KDIGO 2012 clinical practice guideline for the evaluation and  
1135 management of chronic kidney disease. *Kidney International Supplements* **3**,  
1136 1-150 (2013).

1137 39. Levey, A.S., Becker, C. & Inker, L.A. Glomerular filtration rate and  
1138 albuminuria for detection and staging of acute and chronic kidney disease in  
1139 adults: a systematic review. *Jama* **313**, 837-846 (2015).

1140 40. Bikbov, B., *et al.* Global, regional, and national burden of chronic kidney  
1141 disease, 1990–2017: a systematic analysis for the Global Burden of Disease  
1142 Study 2017. *The Lancet* **395**, 709-733 (2020).

1143 41. Liao, Y., Liao, W., Liu, J., Xu, G. & Zeng, R. Assessment of the CKD-EPI  
1144 equation to estimate glomerular filtration rate in adults from a Chinese CKD  
1145 population. *Journal of International Medical Research* **39**, 2273-2280 (2011).

1146 42. Pisano, E.D., *et al.* Contrast limited adaptive histogram equalization image  
1147 processing to improve the detection of simulated spiculations in dense  
1148 mammograms. *Journal of digital imaging* **11**, 193-200 (1998).

1149 43. Liu, P., *et al.* Large-Scale Left and Right Eye Classification in Retinal Images.  
1150 in *Computational Pathology and Ophthalmic Medical Image Analysis* (eds.  
1151 Stoyanov, D., *et al.*) 261-268 (Springer International Publishing, Cham, 2018).

1152 44. He, K., Zhang, X., Ren, S. & Sun, J. *Deep Residual Learning for Image*  
1153 *Recognition*, (2016).

1154 45. Paszke, A., *et al.* *PyTorch: An Imperative Style, High-Performance Deep*  
1155 *Learning Library*, (2019).

1156 46. Kamarudin, A.N., Cox, T. & Kolamunnage-Donà, R. Time-dependent ROC  
1157 curve analysis in medical research: Current methods and applications. *BMC*  
1158 *Medical Research Methodology* **17**(2017).

1159 47. Sundararajan, Mukund, Taly, A. & Yan, Q. Axiomatic Attribution for Deep  
1160 Networks. In *Proceedings of the 34th International Conference on Machine*  
1161 *Learning-Volume 70* **3319**(2017).

1162 48. D, G. Understanding bland altman analysis. *Biochemia medica* **25(2):**  
1163 **141-151**(2015).

1164 49. Breslow, N. & Day, N. Statistical Methods in Cancer Research. Volume  
1165 II--The Design and Analysis of Cohort Studies. *IARC scientific publications*,  
1166 1-406 (1987).

1167

1168

1169 **Figure Legend**

1170

1171 **Figure 1. Schematic illustration of an AI system for detection and incidence**  
1172 **prediction of systemic diseases using retinal fundus images. a,** Model development  
1173 of the AI system. The system made two different types of predictions: “regression  
1174 tasks” for continuous values (eGFR and fasting blood glucose); binary “classification”  
1175 tasks for categorical values (CKD detection and T2DM). The AI prediction is  
1176 generated with an ensemble of model instances. **b,** Application and evaluation of the  
1177 AI system. The left panel: Training on a developmental dataset for assessment of  
1178 CKD/CKD staging/ T2DM. The model was then tested on an independent cohort to  
1179 ensure the generalizability. The middle panel: we developed an AI platform to predict  
1180 future disease development based on fundus images from longitudinal cohorts, which  
1181 should help strengthen CKD/T2DM surveillance. The right panel: a prospective study  
1182 on a “Point of care” setting using a smartphone. Fundus images captured on the phone  
1183 were transmitted to a cloud hosting the AI model which produces an instant report  
1184 transmitted back to the smartphone, an ophthalmologist and a hospital. CKD, chronic  
1185 kidney disease; T2DM, type 2 diabetes mellitus.

1186 **Figure 2. Performance in identifying CKD and early CKD of the AI models.**  
1187 ROC curves using the metadata-only model, the fundus-only model and the combined  
1188 model. **a-c,** ROC curves showing performance of CKD detection on: **a,** internal test  
1189 set (case: control=484:2,685); **b,** external test set 1 (case: control=676:3,880); and **c,**  
1190 external test set 2 (case: control=228:1,014), a prospective point-of-care pilot study  
1191 using retinal fundus images from a smartphone. **d-f,** ROC curves showing  
1192 performance of early CKD detection on: **d,** internal test set (case: control=159:2,685);  
1193 **e,** external test set 1(case: control=240:3,880); **f,** external test set 2 (case:  
1194 control=71:1,014), a prospective point-of-care pilot study using retinal fundus images  
1195 from a smartphone.

1196 **Figure 3. Model performance in assessing GFR from retinal fundus images. a-c,**  
1197 Bland-Altman plot for the agreement between the predicted and estimated GFR. **The**  
1198 X-axis represents the mean of predicted and estimated GFR, and the Y-axis represents  
1199 the difference between the two measurements. Assessing agreement between  
1200 predicted and estimated GFR with mean absolute error (MAE) and intraclass  
1201 correlation coefficient (ICC). **a,** In internal test set. **b,** In external test 1. **c,** In external  
1202 test 2, a prospective point-of-care pilot study.  
1203 **d-f,** Correlation analysis of the predicted eGFR vs actual eGFR generated using the  
1204 regression model, with mean absolute error (MAE), Pearson’s correlation coefficient  
1205 (PCC) and R squared (R2). **d,** AI performance on internal test set. **e,** AI performance  
1206 on the external test set 1. **f,** AI performance on the external test set 2: “point-of-care”  
1207 study.

1208

1209 **Figure 4. Kaplan-Meier plots illustrating the prediction of CKD and advanced+**  
1210 **CKD development.** The y-axis is the survival probability, measuring the probability

of not progressing to a disease outcome. The x-axis is the time in months. Survival curves in different colors are the high-risk, medium-risk and low-risk subgroups stratified by the upper and lower quartiles in the tuning dataset. **a and b**, The incidence of CKD in **a**, the internal longitudinal test set and **b**, in the external longitudinal test set. **c and d**, The incidence of advanced+ CKD (corresponding to stage 3 or more severe) in **c**, internal longitudinal test set and **d**, in the external longitudinal test set. P-value is obtained from the log-rank test.

**Figure 5. Performance in identifying and incidence prediction of T2DM with the AI models.** **a and b**, AI performance in detecting T2DM using the metadata-only model, the fundus-only model and the combined model. **a**, ROC curves showing performance of T2DM detection on internal test set (case: control=2,361:6,366). **b**, ROC curves showing performance of T2DM detection on external test set 1(case: control=2,823:5,236). **c**, ROC curves showing performance of T2DM detection on external test set 2: “point-of-care” study (case: control=1,189:1,892). **d and e**, Kaplan-Meier plots showing the incidence of T2DM during the follow-up visits in **c**, internal longitudinal test set and **d**, external longitudinal test set.

**Figure 6. Gradient visualizations of AI predictions of CKD staging using the Integrated Gradient algorithm.** Visual explanations of the areas of the images most important for the determination of the model prediction for qualitative review and clinical relevance, including **a**, Early CKD (corresponding to stage 1 and 2), **b**, Advanced CKD (corresponding to stage 3), and **c**, Severe+ CKD (corresponding to stage 4+). The columns are (1) the original fundus image, (2) a saliency map, and (3) a saliency map overlaying the original image.

1235 **Table**

1236

1237 **Table 1. Basic characteristics of patients in the developmental dataset and**  
 1238 **external validation cohorts for systemic diseases detection.** The numbers of retinal  
 1239 fundus images used for identifying systemic conditions are shown in each cohort.  
 1240 T2DM, Type 2 diabetes mellitus; CKD, chronic kidney disease; eGFR, estimated  
 1241 glomerular filtration rate; DR, diabetic retinopathy; NDR, diabetes mellitus with no  
 1242 DR; BMI, Body mass index.

1243

Cohorts	Developmental Dataset		Internal test set (CC-FII)	External test set 1	External test set 2: "Point-of-care"
	Training set (CC-FII-C)	Tuning set (CC-FII-C)			
Number of images	60,244	8,614	17,454	16,118	6,162
Number of subjects	30,122	4,307	8,727	8,059	3,081
Male, n (%)	15,325 (50.9%)	2,205 (51.2%)	4,412 (50.6%)	4,441 (55.1%)	1,476 (47.9%)
Age (y), mean (SD)	50.4±14.6	50.6±14.7	50.1±14.6	53.9±13.5	49.0±13.4
BMI (kg/m <sup>2</sup> ), mean (SD)	24.7±2.4	24.7±2.4	24.7±2.5	24.8±2.6	24.6±2.7
Hypertension, n (%)	9,098 (30.2%)	1,308 (30.4%)	2,601 (29.8%)	2,608 (32.4%)	908 (29.5%)
eGFR (mL/min/1.73 m <sup>2</sup> ), mean (SD)	97.1±22.6, n=19,261	97.0±22.3, n=2,773	97.6±21.7, n=5,643	94.7±24.3, n=4,994	98.3±20.6, n=3,065
Blood glucose (mmol/L), mean (SD)	6.3±2.6, n=19,940	6.3±2.7, n=2,884	6.2±2.7, n=5,857	7.0±3.1, n=5,373	6.6±3.0, n=3,039
<b>CKD, n (%)</b>	1,906 (17.4%), n=10,977	251 (16.0%), n=1,569	484 (15.3%), n=3,169	676 (14.8%), n=4,556	228 (18.4%), n=1,242
Early CKD <sup>a</sup> , n (%)	648 (34.0%)	82 (32.7%)	159 (32.9%)	240 (35.5%)	71 (31.1%)
Advanced CKD <sup>b</sup> , n (%)	828 (43.4%)	105 (41.8%)	211 (43.6%)	278 (41.1%)	111 (48.7%)
Severe+ CKD <sup>c</sup> , n (%)	430 (22.6%)	64 (25.5%)	114 (23.6%)	158 (23.4%)	46 (20.2%)
<b>T2DM, n (%)</b>	8,791 (29.2%), n=30,122	1,286 (29.9%), n=4,307	2,361 (27.1%), n=8,727	2,823 (35.0%), n=8,059	1,189 (38.6%), n=3,081
T2DM-DR, n (%)	1,414 (16.1%)	228 (17.7%)	392 (16.6%)	425 (15.1%)	141 (11.9%)
T2DM-NDR, n (%)	7,377 (83.9%)	1,058 (82.3%)	1,969 (83.4%)	2,398 (84.9%)	1,048 (88.1%)
<i>n</i> indicates the number of patients for whom that measurement was available.					

<sup>a</sup>Early CKD is defined as an eGFR of more than 60 mL/min/1.73 m<sup>2</sup> with albuminuria (corresponding to stage 1 and 2).

<sup>b</sup>Advanced CKD is defined as an eGFR of 30-59 mL/min/1.73 m<sup>2</sup> (corresponding to stage 3).

<sup>c</sup>Severe+ CKD is defined as an eGFR of less than 30 mL/min/1.73 m<sup>2</sup> (corresponding to stage 4 and above).

1244

1245

1246

1247

1248

1249

**Table 2. Incidence rates of the CKD/T2DM (per 1000 person-year) on the internal longitudinal test set and the external longitudinal test set according to three-strata of the AI models.**

Cohorts	Prognostic analysis	Risk group	Low risk	Medium risk	High risk	Overall
Internal longitudinal test set (CC-FII-L), n=3,376	CKD	Number of participants	426	827	432	1685
		Number of events	8	23	49	80
		Incidence rate (per 1,000 person-years, 95% CI)	4.4 (1.9, 8.6)	6.5 (4.1, 9.7)	27.0 (20.0, 35.7)	11.1 (8.8, 13.8)
	T2DM	Number of participants	441	868	469	1778
		Number of events	7	25	57	89
		Incidence rate (per 1,000 person-years, 95% CI)	3.9 (1.6, 8.1)	7.4 (4.8, 11.0)	30.0 (22.7, 38.8)	12.6 (10.2, 15.6)
External longitudinal test set, n=2,112	CKD	Number of participants	359	1016	509	1884
		Number of events	13	27	26	66
		Incidence rate (per 1,000 person-years, 95% CI)	9.4 (5.0, 16.0)	7.3 (4.8, 10.6)	14.5 (9.5, 21.3)	9.6 (7.4, 12.2)
	T2DM	Number of participants	370	1370	1404	3144
		Number of events	4	69	118	191
		Incidence rate (per 1,000 person-years, 95% CI)	2.8 (0.7, 7.1)	14.8 (11.5, 18.7)	25.1 (20.7, 30.0)	17.7 (15.2, 20.3)

1250

1251

## Extended Data Tables

**Extended Data Table 1. Characteristics of patients in the developmental set and validation sets of two longitudinal cohorts.** The numbers of retinal fundus images used for predicting the development of systemic conditions are shown in each cohort. T2DM, Type 2 Diabetes Mellitus; CKD, chronic kidney disease; eGFR, estimated glomerular filtration rate; DR, diabetic retinopathy; NDR, diabetes mellitus with no DR.

Longitudinal Cohorts	Developmental Dataset	Internal longitudinal test set (CC-FII-L)	External longitudinal test set
	Training and Tuning set (CC-FII-L)		
Number of images	16,314	4,224	6,752
Number of subjects	8,157	2,112	3,376
Male, n (%)	3,425 (42.0%)	845 (40.0%)	1,426 (42.2%)
Age (y), mean (SD)	46.2±14.4	46.0±14.5	51.8±13.7
BMI (kg/m <sup>2</sup> ), mean (SD)	24.1±3.4	24.1±3.5	24.6±3.6
Hypertension, n (%)	2,518 (30.9%)	649 (30.7%)	1,161 (34.4%)
Follow-up time (months), mean (SD)	51.6±15.8	51.6±15.8	51.1±8.5
<b>Participants with known CKD outcomes</b>	6,467	1,685	1,884
Diabetes, n(%)	1,688 (26.1%)	414 (24.6%)	456 (24.2%)
CKD outcome (to Early CKD)	160 (2.5%)	39 (2.3%)	50 (2.7%)
CKD outcome (to Advanced+ CKD)	148 (2.3%)	41 (2.4%)	16 (0.8%)
<b>Participants with known T2DM outcomes</b>	6,807	1,778	3,144
T2DM outcome (to T2DM)	396 (5.8%)	89 (5.0%)	191 (6.1%)

**Extended Data Table 2. AI Performance for detection of CKD or T2DM using logistic regression models on internal and external test sets.**



1264

Cohorts	Internal test set (CC-FII)	External test set 1	External test set 2: “Point-of-care”
CKD (LR)	0.814 (0.795-0.830)	0.801 (0.785-0.816)	0.784 (0.751-0.813)
T2DM (LR)	0.773 (0.756-0.786)	0.788 (0.767-0.806)	0.774 (0.740-0.802)

1265

1266

1267

1268

1269

**Extended Data Table 3. Univariate and multivariate survival analyses of CKD/T2DM conducted using Cox proportional hazards methods (likelihood ratio test).**

Prognostic analysis	Covariates	Univariate analysis		Multivariate analysis	
		Hazard ratio	p-value	Hazard ratio	p-value
CKD	Age	1.05 (1.05-1.06)	<0.001	1.03 (1.02-1.03)	<0.001
	Sex	0.76 (0.63-0.91)	0.004	0.73 (0.58-0.92)	0.007
	Hypertension	3.53 (2.90-4.30)	<0.001	1.54 (1.23-1.93)	<0.001
	BMI	1.08 (1.05-1.11)	<0.001	1.03 (0.99-1.08)	0.112
	Height	0.97 (0.96-0.98)	<0.001	0.98 (0.97-1.00)	0.059
	Weight	1.01 (1.00-1.01)	0.149	1.00 (0.99-1.02)	0.840
	Diabetes	5.12 (4.17-6.28)	<0.001	2.51 (2.00-3.14)	<0.001
	Fundus (per standard deviation)	3.94 (3.46-4.49)	<0.001	2.10 (1.77-2.50)	<0.001
T2DM	Age	1.04 (1.03-1.04)	<0.001	1.02 (1.02-1.03)	<0.001
	Sex	0.65 (0.56-0.76)	<0.001	0.98 (0.79-1.20)	0.823
	Hypertensions	3.24 (2.76-3.79)	<0.001	1.62 (1.36-1.93)	<0.001
	BMI	1.16 (1.14-1.18)	<0.001	1.09 (1.05-1.12)	<0.001
	Height	1.00 (0.99-1.01)	0.807	0.99 (0.98-1.01)	0.425
	Weight	1.03 (1.03-1.04)	<0.001	1.02 (1.01-1.03)	0.002
	Fundus (per standard deviation)	4.35 (3.64-5.19)	<0.001	1.78 (1.37-2.32)	<0.001

1270

1271

1272

1273

1274

1275

1276

1277

1278

**Extended Data Table 4. Performance of progression prediction model to CKD or advanced+ CKD event based on the metadata-only model, and the combined model (including fundus images and metadata) on the internal and external test sets.** Concordance index (C-index) for right-censored data and 95% confidence intervals (CI) measure the model performance by comparing the progression information (disease labels and progression days) with predicted risk scores. A larger C-index correlates with better progression prediction performance.

Tasks	Progression prediction models	Internal longitudinal test set	External longitudinal test set
CKD	Combined model	0.845 (95% CI: 0.789-0.910)	0.719 (95% CI: 0.627-0.807)
	Metadata model	0.756 (95% CI: 0.699-0.810)	0.651 (95% CI: 0.569-0.730)
Advanced+CKD	Combined model	0.933 (95% CI: 0.909-0.955)	0.912 (95% CI: 0.823-0.972)
	Metadata model	0.847 (95% CI: 0.804-0.896)	0.832 (95% CI: 0.720-0.924)

**Extended Data Table 5. Performance of progression prediction model to T2DM event based on the metadata-only model, and the combined model (including fundus images and metadata) on the internal and external test sets.** Concordance index (C-index) for right-censored data and 95% confidence intervals (CI) measure the model performance by comparing the progression information (disease labels and progression days) with predicted risk scores. A larger C-index correlates with better progression prediction performance.

Tasks	Progression prediction models	Internal longitudinal test set	External longitudinal test set
T2DM	Combined model	0.781 (95% CI: 0.743-0.819)	0.765 (95% CI: 0.723-0.799)
	Metadata model	0.774 (95% CI: 0.732-0.819)	0.746 (95% CI: 0.706-0.775)

**Extended Data Table 6. Performance of the AI system for CKD detection (including Early CKD, Advanced and Severe+) from normal controls using retinal fundus images.** Each row represents metrics based on the corresponding operation point set to perform with high NPV and PPV for CKD screening. CI, confidence interval; PPV, positive predictive value; NPV, negative predictive value.

Operating point based on the tuning set	Cohorts	Sensitivity	Specificity	Reliability of computer-aided decision (CAD)
Positive	Internal test set	43.3% (38.9-49.4)	99.4% (99.1-99.7)	PPV: 92.4% (88.3-95.7)
	External test set 1	34.8% (31.8-38.6)	99.2% (98.9-99.5)	PPV: 88.4% (83.9-92.8)
	External test set 2: "Point-of-care"	29.9% (21.5-34.3)	99.2% (98.5-99.7)	PPV: 89.3% (80.8-95.5)
Negative	Internal test set	99.3% (98.2-100.0)	42.8% (41.0-44.5)	NPV: 99.7% (99.3-100.0)
	External test set 1	99.4% (98.8-99.8)	37.5% (36.1-38.9)	NPV: 99.7% (99.5-99.9)
	External test set 2: "Point-of-care"	99.1% (97.7-100.0)	32.1% (29.0-35.5)	NPV: 99.4% (98.5-100.0)

**Extended Data Table 7. Performance of the AI system for T2DM detection using retinal fundus images.** Each row represents metrics based on the corresponding operation point set to perform with high NPV and PPV for T2DM screening. CI, confidence interval; PPV, positive predictive value; NPV, negative predictive value.

Operating point based on the tuning set	Cohorts	Sensitivity	Specificity	Reliability of computer-aided decision (CAD)
Positive	Internal test set	59.1% (54.8-62.0)	97.8% (97.4-98.1)	PPV: 78.7% (75.1-82.1)
	External test set 1	15.9% (12.6-20.7)	99.6% (99.5-99.8)	PPV: 77.9% (71.1-86.2)
	External test set 2: "Point-of-care"	12.1% (8.5-16.1)	99.5% (99.2-99.8)	PPV: 72.7% (59.4-86.1)
Negative	Internal test set	99.3% (98.5-99.8)	41.6% (40.4-43.0)	NPV: 99.8% (99.5-99.9)
	External test set 1	98.8% (97.8-99.7)	31.9% (30.9-33.3)	NPV: 99.7% (99.4-99.9)
	External test set 2: "Point-of-care"	98.5% (96.5-100.0)	44.4% (42.0-46.5)	NPV: 99.6% (99.1-100.0)

**Extended Data Table 8. Basic characteristics of patients in the Multi-ethnicity validation cohort for systemic diseases detection.** Shown are the numbers of retinal fundus images used for identifying systemic conditions. T2DM, Type 2 diabetes mellitus; CKD, chronic kidney disease; eGFR, estimated glomerular filtration rate; DR, diabetic retinopathy; NDR, diabetes mellitus with no DR; BMI, Body mass index.

Cohort	Multi-ethnicity validation set
Number of images	1,230
Number of subjects	615
Male, n (%)	304 (49.4%)
Age (y), mean (SD)	60.5±12.9
BMI (kg/m <sup>2</sup> ), mean (SD)	25.3±3.1
Hypertension, n (%)	260 (42.3%)
eGFR (mL/min/1.73 m <sup>2</sup> ), mean (SD)	86.2±19.0, n=577
Blood glucose (mmol/L), mean (SD)	7.0±2.7, n=586
CKD, n (%)	93 (21.2%), n=439
T2DM, n (%)	343 (55.8%), n=615

## 1314 **Extended Data Figure Legend**

1315

1316 **Extended Data Figure 1. The flowchart of the AI platform with an ensemble of**  
1317 **model instances.** We first developed retinal fundus image enhancement models using  
1318 color normalization and contrast-limited adaptive histogram equalization (CLAHE)  
1319 techniques. Four types of fundus images after the application of color normalization  
1320 and CLAHE image enhancements: original image, image after applying the CLAHE  
1321 transformation only, image after applying the color normalization transformation only,  
1322 and image after applying both the CLAHE and color normalization transformations.  
1323 Each image instance separately makes a prediction, and these are combined by  
1324 averaging the results for producing a robust AI model.

1325

1326 **Extended Data Figure 2. Flow diagram describing the datasets used for our AI**  
1327 **system for CKD/T2DM detection and incidence prediction.** Patient inclusion and  
1328 exclusion criteria were also considered.

1329

1330 **Extended Data Figure 3. Model performance in assessing GFR/CKD staging**  
1331 **using retinal fundus images. a-c,** Bland-Altman plot for predicted and actual eGFR  
1332 after calibrating the model output. AI performance on **a**, the internal test set, **b**, the  
1333 external test set 1 and **c**, the external test set 2: “point-of-care” study. **d and e**, AI  
1334 performance in detecting severe+ CKD from other stages of CKD (early and  
1335 advanced CKD) with the “regression model” and “classification model”. The blue  
1336 curve denoted “classification model” using retinal fundus images. The orange curve  
1337 denoted “regression model” using thresholds of the predicted GFR from retinal  
1338 fundus images. **d**, In the internal test. **e**, In the external test set 1. **f**, Correlation  
1339 analysis of the predicted eGFR of the right eye versus the predicted eGFR of the left  
1340 eye in normal, early CKD (stages 1 and 2), and CKD. ICC, intraclass correlation  
1341 coefficient.

1342

1343 **Extended Data Figure 4. Prediction of fasting blood glucose using retinal fundus**  
1344 **images. a-c,** Correlation analysis of predicted blood glucose vs actual blood glucose  
1345 generated using the regression model. **a**, the internal test set, **b**, the external test set 1  
1346 and **c**, the external test set 2: “point-of-care” study. **d-f**, Bland-Altman plot for the  
1347 agreement between the predicted and actual blood glucose levels (mmol/L). Assessing  
1348 agreement between predicted and estimated blood glucose with mean absolute error  
1349 (MAE) and intraclass correlation coefficient (ICC). **d**, the internal test set. **e**, the  
1350 external test set 1. **f**, the external test set 2: the prospective ‘point-of care’ pilot study.  
1351 **g-i**, Bland-Altman plot for predicted and actual blood glucose after calibrating the  
1352 model output. AI performance on **g**, internal test set, **h**, the external test set 1 and **i**,  
1353 the external test set 2: the ‘point-of-care’ study.

1354

1355 **Extended Data Figure 5. AI performance at detecting T2DM patients with**  
1356 **images with no apparent signs of diabetic retinopathy (NDR) and images with**  
1357 **diabetic retinopathy (DR). a and b**, Visualizations for detecting T2DM patients by

1358 highlighting the regions the AI model focuses. **a**, T2DM, with “no signs of diabetic  
1359 retinopathy” (NDR), and **b**, T2DM, with severe DR. The rows are (1) the original  
1360 fundus image, (2) a saliency map, and (3) a saliency map overlaying the original  
1361 image. **c**, ROC curves showing performance of binary classification models in the  
1362 internal test set. Comparison of the AI performance at detecting of T2DM patients  
1363 with only images with NDR compared to the performance with both images with DR.

1364 **Extended Data Figure 6. Performance of the AI model on the external**  
1365 **multi-ethnicity validation cohort from Kashi and Macau. a and b, ROC curves**  
1366 **showing performance of the metadata-only model, the fundus-only model and the**  
1367 **combined model on the classification of systematic diseases: a, CKD and b, T2DM.**

1368 **Figure 2. Performance in identifying CKD and early CKD of the AI models.**  
1369 ROC curves using the metadata-only model, the fundus-only model and the combined  
1370 model. **a-c**, ROC curves showing performance of CKD detection on: **a**, internal test  
1371 set (case: control=484:2,685); **b**, external test set 1 (case: control=676:3,880); and **c**,  
1372 external test set 2 (case: control=228:1,014), a prospective point-of-care pilot study  
1373 using retinal fundus images from a smartphone. **d-f**, ROC curves showing  
1374 performance of early CKD detection on: **d**, internal test set (case: control=159:2,685);  
1375 **e**, external test set 1 (case: control=240:3,880); **f**, external test set 2 (case:  
1376 control=71:1,014), a prospective point-of-care pilot study using retinal fundus images  
1377 from a smartphone.

1378 **Extended Data Figure 7. Prediction of the development of CKD and T2DM using**  
1379 **the metadata-only model. Kaplan Meier plot illustrating the incidence of**  
1380 **CKD/T2DM.** The blue, orange, and green lines represent stratified scores for low risk,  
1381 medium risk, and high risk, respectively, using the upper and lower quartiles in the  
1382 tuning dataset. **a and b**, progression to CKD on **a**, internal longitudinal test set, **b**,  
1383 external longitudinal test set. **c and d**, Progression to advanced+ CKD on **c**, internal  
1384 longitudinal test set, **b**, external longitudinal test set. **e and f**, Progression to T2DM on,  
1385 **e**, internal longitudinal test set, **f**, external longitudinal test set.

1386  
1387 **Extended Data Figure 8. The cumulative hazard functions of three stratified risk**  
1388 **subgroups (high- medium- and low-risk) using the combined progression**  
1389 **prediction model** on (a) the internal longitudinal test set and (b) the external  
1390 longitudinal test set. The solid line is the mean cumulative hazard scores at each time  
1391 point. The area of the same color represents the 95% confidence interval. Numbers at  
1392 risk are also included.

1393  
1394 **Extended Data Figure 9. Prediction of the development of CKD and T2DM using**  
1395 **Time-dependent ROC curves. a and b, ROC curves for quantifying AI model**  
1396 **performance for the incidence of CKD: a**, in the internal longitudinal test set for 5  
1397 years follow up (case: control=62:470). **b**, in the external longitudinal test set for 4  
1398 years follow up (case: control=40:663). **c and d**, ROC curves for quantifying AI  
1399 model performance for the incidence of T2DM. **c**, in the internal longitudinal test set

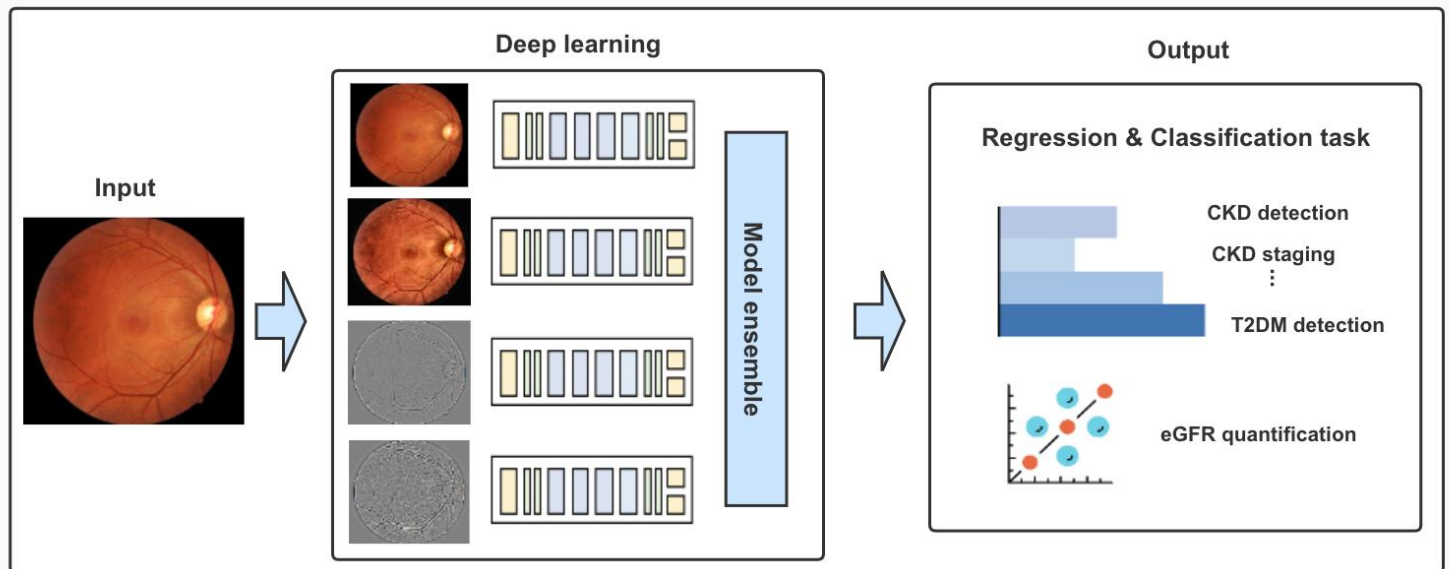
1400 for 5 years follow up (case: control=68:425). **d**, in the external longitudinal test set for  
1401 4 years follow up (case: control=96:1,266).

1402

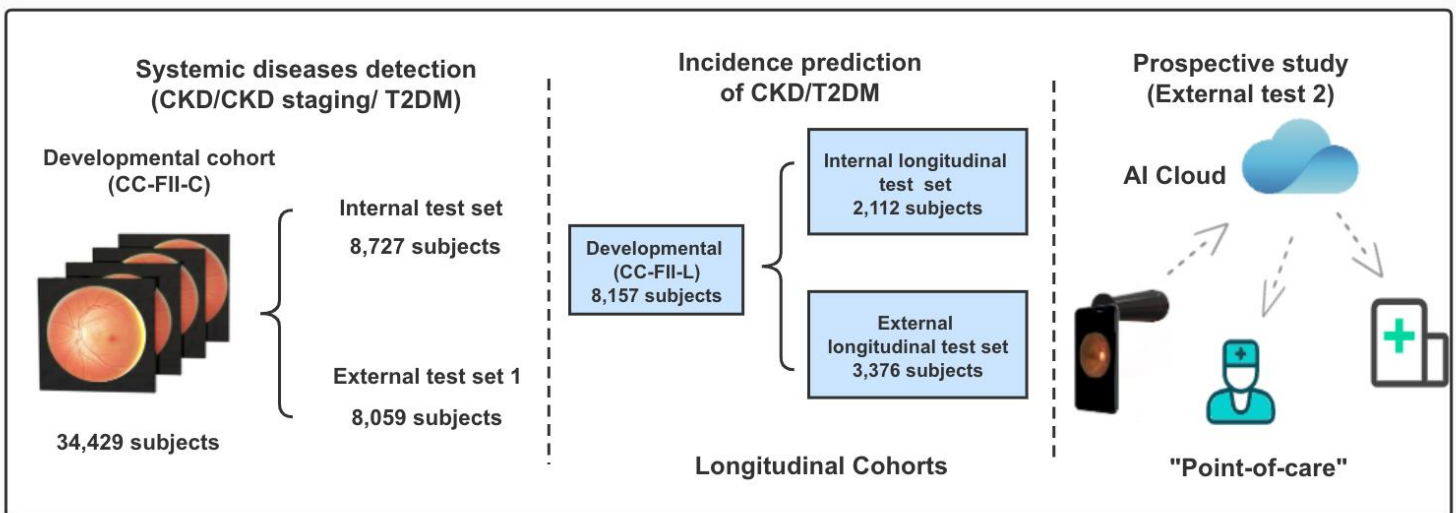
1403 **Extended Data Figure 10. Design Illustration on hand-held smartphone camera**  
1404 **attachment.** We used a standard 3D printer to make a customized adaptor that can be  
1405 fitted and attached to an iPhone X.

**Figure 1**

**a**



**b**



**Figure 2**

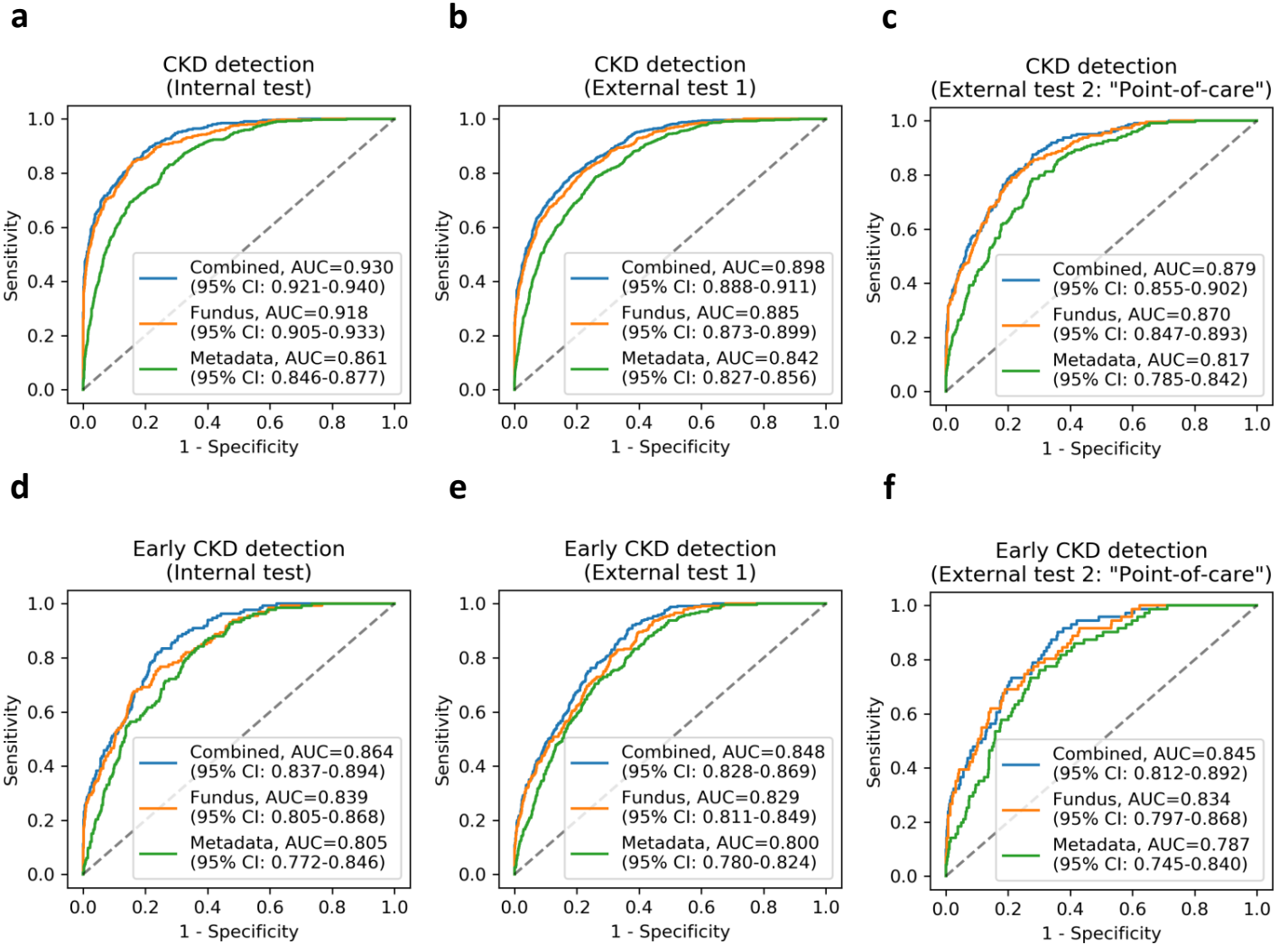




Figure 3

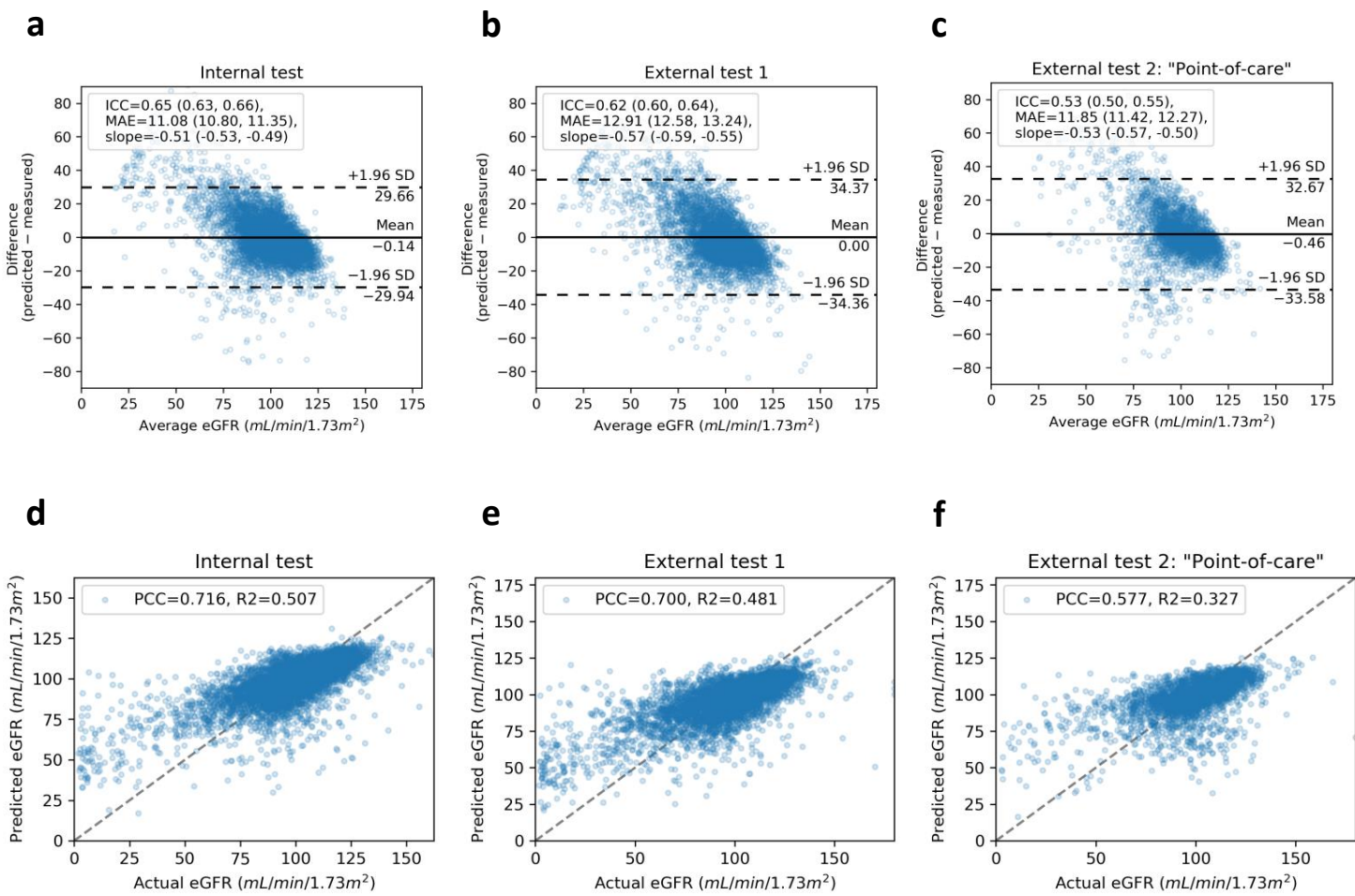
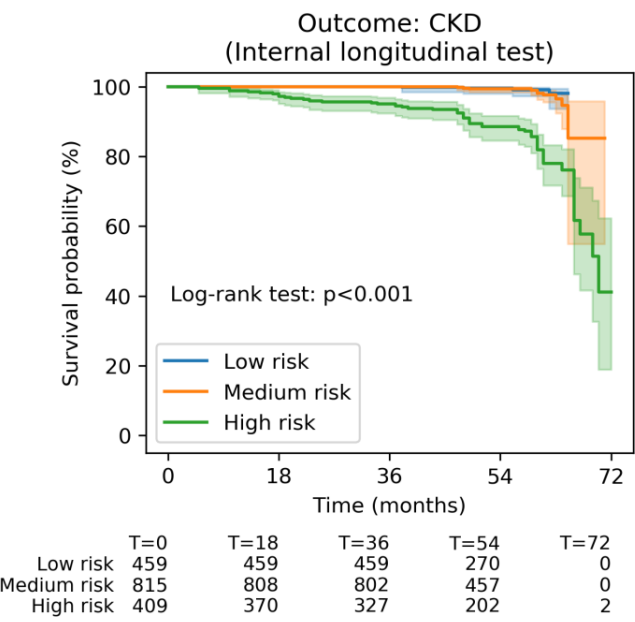
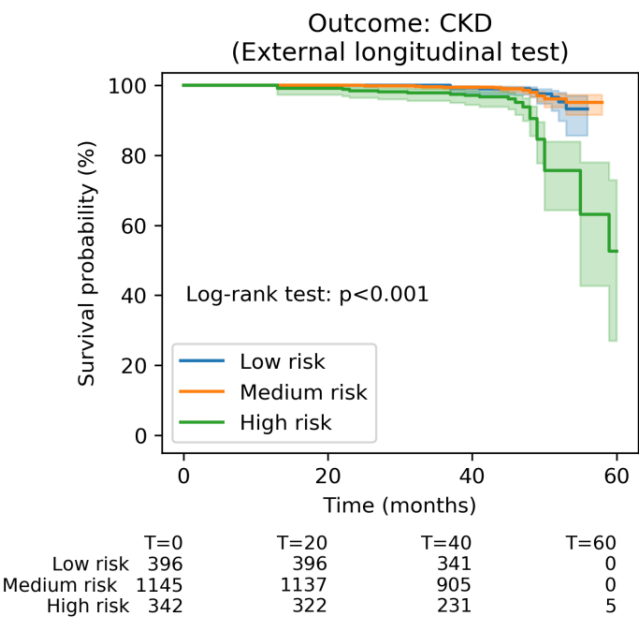


Figure 4

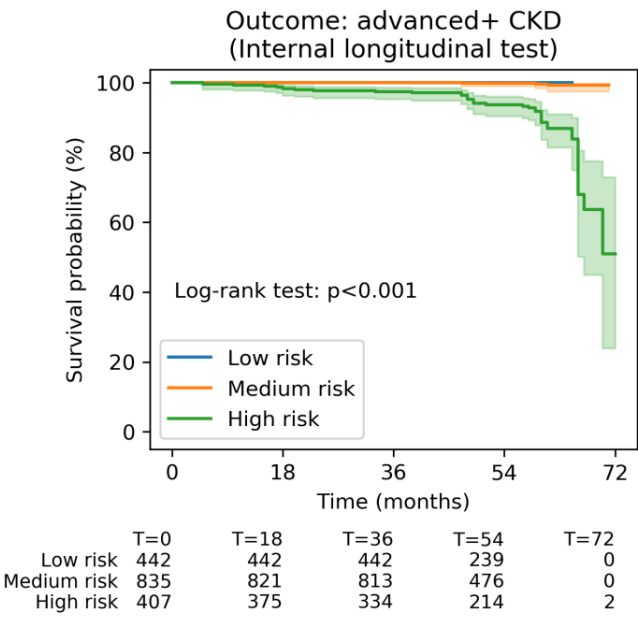
a



b



c



d

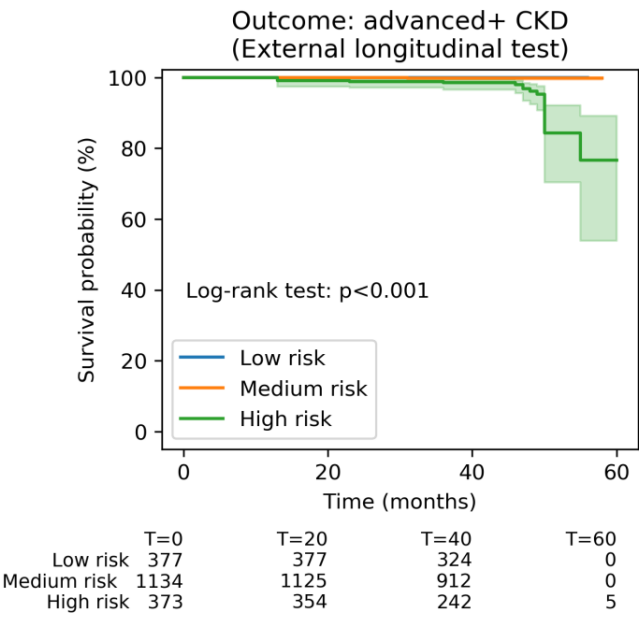


Figure 5

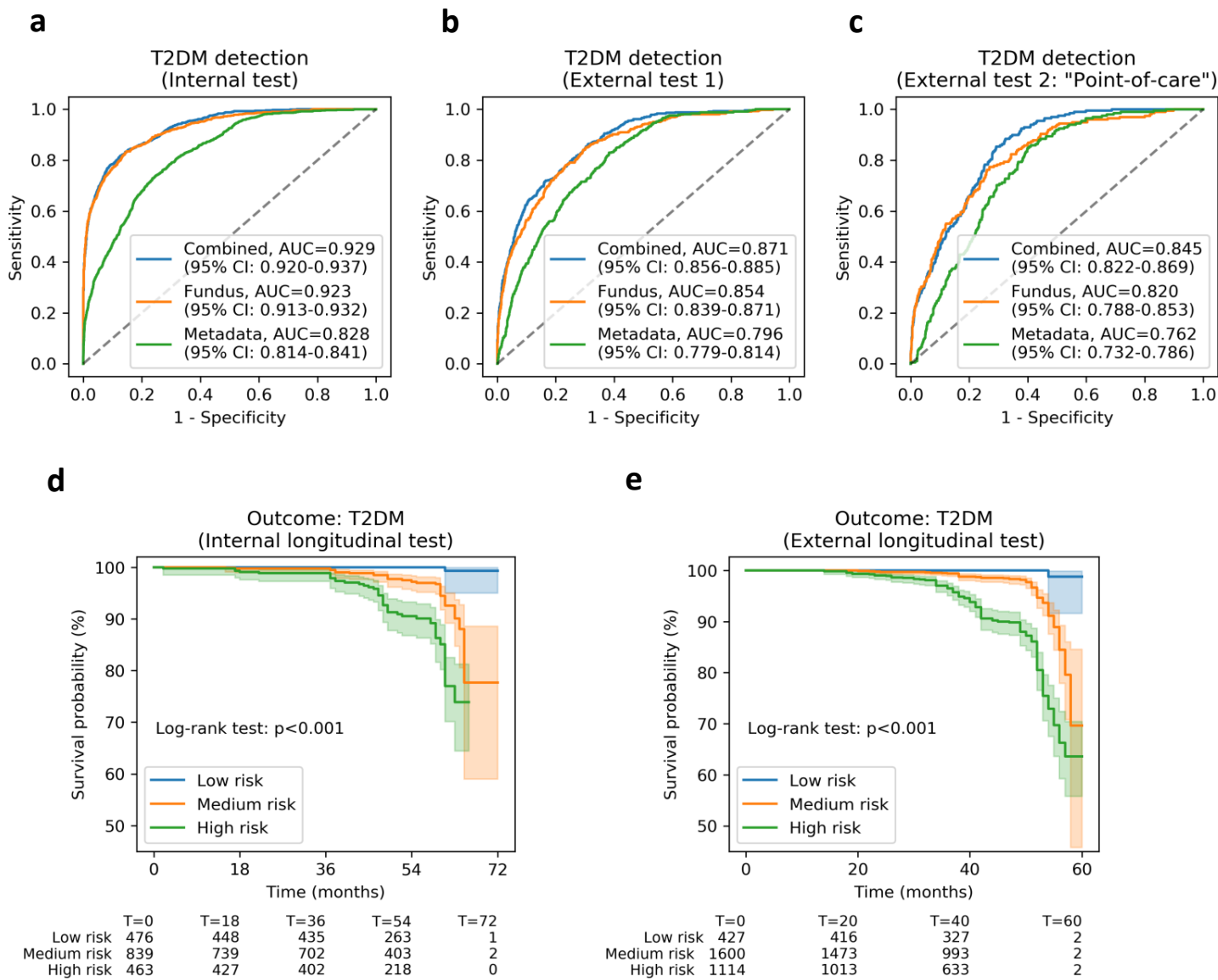


Figure 6

