

Statistical machine learning of sleep and physical activity phenotypes from sensor data in 96,609 UK Biobank participants.

Matthew Willetts^{1†}, Sven Hollowell^{2,3†}, Louis Aslett¹, Chris Holmes^{1,2‡}, Aiden Doherty^{2,3,4‡*}

1 Department of Statistics, University of Oxford, Oxford, U.K. **2** Big Data Institute, Li Ka Shing Centre for Health Information and Discovery, University of Oxford **3** Nuffield Department of Population Health, BHF Centre of Research Excellence, University of Oxford, Oxford, U.K. **4** Institute of Biomedical Engineering, Department of Engineering Science, University of Oxford, Oxford, U.K.

* Corresponding author. E-mail: aiden.doherty@bdi.ox.ac.uk

† Joint first authors.

‡ Joint last authors.

ABSTRACT

Current public health guidelines on physical activity and sleep duration are limited by a reliance on subjective self-reported evidence. Using data from simple wrist-worn activity monitors, we developed a tailored machine learning model, using balanced random forests with Markov confusion matrices, to reliably detect a number of activity modes. We show that physical activity and sleep behaviours can be classified with 87% accuracy in 84,616 minutes of recorded free-living behaviours from 57 adults. These trained models can be used to infer fine resolution activity patterns at the population scale in 96,609 participants. For example, we find that men spend more time in both low- and high- intensity behaviours, while women spend more time in mixed behaviours. Walking time is highest in spring and sleep time lowest during the summer. This work opens the possibility of future public health guidelines informed by the health consequences associated with specific, objectively measured, physical activity and sleep behaviours.

INTRODUCTION

The way that adults spend their time (for example how long they sleep, walk, and sit) has important health implications¹⁻⁵. However this evidence is largely based on self-reported data that are crude and prone to measurement error⁶. Therefore, uncertainty exists on the exact amount and types of sleep and physical activity behaviours that should be recommended, and which interventions and programmes may be most effective in helping people live more healthily. As a result, longitudinal studies such as UK Biobank^{7,8} now aim to collect objective measures of sleep and physical activity via wrist-worn accelerometers so that their health consequences can be understood. As an important first step, 'vector magnitude' methods have been developed to objectively measure the volume and intensity levels of physical activity from accelerometer data in large health datasets⁹⁻¹¹. However, a better understanding of the health consequences of individual lifestyle health behaviours (such as sleeping, sitting, and walking) would arguably help inform public health recommendations that are readily interpretable and implementable. For example, a recommendation of 30 min/day of walking might be more easily understood and actionable than 30 min/day of moderate-to-vigorous intensity physical activity¹².

Flaws exist in the validation of current methods to extract behavioural information from accelerometer data for relevant biomedical analysis. For example, machine learning methods to detect specific behaviours of interest¹³, such as walking and sitting, have generally not been validated in realistic free-living environments¹⁴. The validation of these methods in laboratory scenarios is unrealistic as it usually involves a limited number of activities¹⁵, poor variety within each activity, and an unrealistic relative contribution in time for each activity type¹⁶. As a result, it is difficult to verify whether current assumptions on the ability to predict walking¹⁷, or bicycling^{14,18}, from sensor data are correct or biased in some manner. Recent work by Ellis and colleagues in a US study has demonstrated that the relevance and accuracy of accelerometer based machine learning methods improves across a range of activities when trained on free-living, rather than controlled laboratory data¹⁹. However, machine learning methods have not been assessed in large scale health sensor datasets for face validity or investigated to evaluate if they offer behavioural insight.

In this paper we describe the development of a machine learning method to objectively measure lifestyle health behaviours from wrist-worn accelerometer data. We firstly assessed its performance in free-living scenarios using a dataset of 57 adults, 40 of which are female, aged 18-91 who wore an accelerometer and wearable camera (a method comparable to direct-observation²⁰). This labelled dataset for machine learning development and held-out validation is many times larger (~84,616 minutes of behaviour) than previous lab-based studies [330 - 3,600 mins^{14,17,21}] and free-living studies with short periods of direct observation [3,400 - 24,000 mins^{22,23}]. We then report the utility of our trained method to assess behavioural variation in more than 100,000 UK Biobank participants aged 43-78 by different self-reported phenotypes. This approach provides an automated analysis of objectively measured behavioural variation in lifestyle behaviours and can be used by researchers to study social and health behaviours at a resolution not previously available.

RESULTS

*** Table 1 around here ***

For activity recognition, we trained a balanced random forest with a Markov confusion matrix containing transitions between predicted activity states and emissions trained using a free-living groundtruth to identify six pre-defined classes of behaviour {bicycling, sit/stand, walking, vehicle, mixed activity, sleep} from accelerometer data. Full details of these models are provided in MATERIALS and METHODS, subsections on activity recognition and time smoothing. For comparison against a free-living groundtruth and to maximise available training data, we conducted leave-one-subject-out cross validation for each of our 57 participants. Over these our model obtained a mean accuracy of 86% with a kappa inter-rater agreement score of 0.79 over all the behaviour types in 30-second windows. Sleep/wake classification was most robust, see our minute-level confusion matrix (Table 1). As expected, there was a wide range of individual variation in classification performance at the daily level, see Bland Altman plots for each activity type (Fig. S1). Overall classification performance was not materially altered by the inclusion of age and sex as parameters (Fig. S2), or increasing the number of decision trees in the random forest (Fig. S3). However, the inclusion of hidden Markov model time smoothing, over base random forest predictions, boosted overall classification performance Kappa score from 0.67 to 0.79 (Table S1). For energy expenditure prediction, we pre-specified 9 classes of behaviour {sleep, sit/stand+activity, sit/stand+low-activity, walking+activity, sit/stand, walking, vehicle, bicycling, sports} that were based on grouping scores in Metabolic Equivalent of Task²⁴ (MET). We then took the marginal probability of each state and used this to create a weighted average of the MET scores from the 9 classes. Performing leave-one-subject-out cross validation, we find our model had a root mean squared error of 1.81 MET hours/day ($r=0.87$). This compares favourably (RMSE of 2.34 MET hours/day and $r=0.81$) to using random forests for regression²⁵ on our dataset.

*** Figure 1 around here ***

To assess the face validity of our activity recognition method in a large prospective dataset, we applied our model to 103,712 UK Biobank participants. We removed 7,103 participants who did not wear the device for a sufficient amount of time, or who had device errors²⁶. On 96,609 participants, we then plotted each aggregated activity by time-of-day and stratified groups by characteristics self-reported during the study baseline visit a mean of 5.7 years before accelerometer wear (see Fig. 1). Fig. 1a shows self-reported 'evening' people were more likely to be classified as sleeping at 8am on weekends vs. 'morning' people (44% vs. 25%, $p<10^{-100}$ after adjustment for age, sex, ethnicity, area deprivation, smoking, alcohol, and self-rated overall health). Self-reported car users were more likely to be classified as driving at 8am on weekdays (7.7% vs. 3.9%, $p<10^{-100}$ after adjustment for other factors) (see Fig. 1b). Similarly, self-reported cyclists were more likely to be classified as cycling at 8am on weekdays (3.7% vs. 0.6%, $p<10^{-100}$ after adjustment for other factors). Those in active occupations were more likely not to be classified as sitting or standing at 11am on weekdays than office based workers (50% vs. 32%, $p<10^{-100}$ after adjustment for other factors) (Fig. 1d). Older adults (aged 65+) were more likely to be classified as walking at 11am on weekdays than younger adults (aged <55) (15.8% vs. 11.5%, $p=2.4\times 10^{-24}$ after adjustment)

(Fig. 1e). Finally, retired people were more likely to be classified as doing mixed activity at 11am on weekdays than their working counterparts (28% vs. 23%, $p=7.0 \times 10^{-32}$ after adjustment). As expected, these age/occupation differences mostly disappear on weekends (Fig. 1d-f).

*** Table 2 around here ***

Table 2 describes the variation in accelerometer-measured total time for each behaviour type, by age, self-rated health, time-of-day, weekday/weekend, and season. Younger participants spent more time in active behaviours than their older counterparts (e.g. 7.3% vs. 6.1% time for walking, $p < 10^{-100}$). However, these age differences for specific behaviours weren't as pronounced as for vector magnitude, which is a proxy measure of overall physical activity⁹. The behaviour patterns of men appeared more polarised than that of women, with more time in low-intensity activities such as sitting (32.5% vs. 35.3%, $p < 10^{-100}$) but also more time in purposeful activity behaviours such as walking (7.2% vs. 6.2%, $p < 10^{-100}$) and bicycling. Women spent more time engaged in mixed activity behaviours than men (17.4% vs. 13.2%, $p < 10^{-100}$). Self-rated health differences were strongest for traditional vector magnitude measures, but also noticeable for sitting time between those in excellent versus fair/poor self-rated health (33.0% vs. 35.3%, $p < 10^{-100}$). While overall physical activity differences between weekdays and weekends were small (Cohen's $d=0.05$), behavioural differences were more noticeable ($d=0.19$ and $d=0.22$ for longer sleep and less sitting time at weekends respectively). Small seasonal differences also existed with walking time highest in spring, and in summer than in winter there was less sleep (34.1% vs. 35.2%, $p=2 \times 10^{-83}$), more bicycling (0.6% vs. 0.4%, $p < 10^{-100}$), and higher energy expenditure (39.6 vs. 38.9 MET hours/day, $p < 10^{-100}$).

*** Fig. 2 around here ***

Due to the time-series nature of the data, it is possible to illustrate the likelihood for each activity state throughout the day in 96,609 UK Biobank participants (see Fig. 2), and how this relates to daily energy expenditure (see Fig. S5). Apart from mixed-activity and METs ($r=0.77$), the different activity types are weakly correlated (overall mean $r=0.18$), indicating their utility as new sources of information (Fig. S4). To support other researchers test or generate new hypotheses using our data, we have created a tool to show behavioural variation by 40+ participant characteristics (see <http://gas.ndph.ox.ac.uk/aidend/activityClassificationPaper/>)

DISCUSSION

This study represents the largest ever assessment of objectively measured sleep and physical activity behaviours using state-of-the-art machine learning methods trained with a free-living groundtruth. To our knowledge, this is the first time that {"bicycling", "mixed", "sit/stand", "sleep", "vehicle", "walking"} behaviours have been objectively measured in a large-scale dataset. We have demonstrated the feasibility of our method as scale in 96,609 UK Biobank participants where, for example, commute times can be seen for those who self-report as cycling to work. The objective and fine-grained measures (with ~20,000

behavioural predictions per person per week) that we have developed will help more precisely understand the effectiveness of treatments and also the disease processes associated with behaviour variation.

The overall classification score of kappa=0.79 for our method (a balanced random-forest with Markov transitions on predicted states and emissions trained in a naturalistic free-living scenario) represents a substantial level of agreement with the wearable camera groundtruth²⁷. This is comparable to the level of performance (kappa=0.80) that we expect from humans annotating behaviour from wearable camera data^{28,29}. Our overall crude accuracy of 86% is at least as good as reported in other free-living studies with wrist-worn accelerometers [85% in³⁰ and 61% in²²]. Direct comparison across such studies is difficult due to heterogeneity across populations, devices, and definition of behaviour labels. For energy expenditure prediction, random forests for regression²⁵ perform better than linear models for wrist-worn data²¹, and we found the inclusion of activity labels reduced noise and thus further improved performance (e.g. movement is high when in a vehicle, but energy expenditure is low). This mirrors findings from studies that used older hip-worn accelerometers³¹.

Our study shows that bicycling can be reliably detected from accelerometer data, an activity previously difficult to classify in laboratory studies^{14,18}. Potential explanations for this might be our use of devices with higher sampling rates (100Hz vs. 1Hz)¹⁸ that can capture important bicycling activities, or models developed in free-living moving bicycles rather than stationary laboratory bicycles¹⁴. Previous laboratory studies indicated that walking can be classified with a high degree of accuracy (sensitivity and specificity both >90%). However, our data and that from Ellis et al.³⁰ shows that walking is challenging to classify in free-living conditions (sensitivity=0.51, specificity=0.97), probably due to it being part of many everyday activities.

While the participants in our free-living dataset are not a random subsample of the UK Biobank study, the size of our training dataset has helped provide a diverse set of representative behaviours, rather than individuals, which is important for model development. The only other study comparable in training set size used a similar study procedure³⁰, but in a US population subgroup and thus could not be reliably extrapolated to UK Biobank data. Our study also uses more stringent evaluation criteria (kappa versus balanced accuracy) that consider unbalanced free-living data where infrequent behaviours are more susceptible to misclassification. Our method is device agnostic, and could be reused in other large sensor datasets^{9,10,32,33}, provided model tuning takes place in a relevant population with free-living groundtruth validation tools such as wearable cameras^{12,30}. For this study we did not use traditional cumbersome methods to collect sleep³⁴ and energy expenditure³⁵ groundtruth data, as we preferred to use proxy reference methods for free-living assessment at scale^{35,36}. We have not generalised the overall descriptive findings to the UK population since the UK Biobank was established as an aetiological study rather than one aimed at population surveillance^{7,8}.

In summary, we describe the first application, to our knowledge, of machine learning to objectively measure lifestyle health behaviours from sensor data in a large prospective health study. Our method has demonstrated substantial agreement with a free-living groundtruth, and shows face validity in a large health dataset. It is now possible to study the

sociological and health consequences of behaviour variation in unprecedented detail. The summary variables that we have constructed are now part of the UK Biobank dataset and can be used by researchers as exposures, confounding factors or outcome variables in future health analyses.

MATERIALS and METHODS

Participants

For the development and free-living evaluation of accelerometer machine learning methods, we used the first 57 participants recruited to the CAPTURE-24 study where adults aged 18-91 were recruited from the Oxford region in 2014-2015³⁷. Participants were asked to wear a wrist-worn accelerometer for a 24-hour period and then given a £20 voucher for taking part in this study that received ethical approval from University of Oxford (Inter-Divisional Research Ethics Committee (IDREC) reference number: SSD/CUREC1A/13-262). For extrapolation to a large health dataset, we used the UK Biobank dataset where 103,712 participants agreed to wear a wrist-worn accelerometer for a seven day period between 2013-2015⁹. UK Biobank is a large prospective study of 500,000 participants that has collected, and continues to collect, extensive phenotypic and genotypic details about its participants, with ongoing longitudinal follow-up for a wide range of health-related outcomes⁷. Demographic and behavioural variables were recorded by a self-completed touchscreen questionnaire during clinic visits between 2006-2010 (see appendix 1). This study (UK Biobank project #9126) was covered by the general ethical approval for UK Biobank studies from the NHS National Research Ethics Service on 17th June 2011 (Ref 11/NW/0382).

Accelerometer

Participants in both studies were asked to wear an Axivity AX3 wrist-worn triaxial accelerometer on their dominant hand at all times. It was set to capture tri-axial acceleration data at 100Hz with a dynamic range of +-8g. This device has demonstrated equivalent signal vector magnitude output on multi-axis shaking tests³⁸ to the GENEActiv accelerometer which has been validated using both standard laboratory and free-living energy expenditure assessment methods^{36,39}.

Groundtruth

To construct a groundtruth of reference behaviours, participants in the Oxford study were asked to wear a Vicon Autographer wearable camera while awake on the study measurement day. Wearable cameras automatically take photographs every ~20 seconds, have up to 16 hours battery life and storage capacity for over one week's worth of images⁴⁰. When worn, the camera is reasonably close to the wearer's eye line and has a wide-angle lens to capture everything within the wearer's view⁴¹. Each image is time-stamped so duration of active travel⁴², sedentary behaviour²⁹, and a range of other physical activity behaviours⁴³ can be captured. Camera data strongly agrees with more expensive direct observation methods to classify activity types [$\kappa=0.92^{20}$]. We used specific ethical guidance for wearable camera research to inform the development of protocols⁴⁴. Images were annotated by human annotators using codes from the compendium of physical activities²⁴, using specific wearable camera browsing software⁴⁵ (Doc S1). For quality

control, our annotators firstly had to achieve a kappa inter-rater agreement score of >0.8 on separate training data. To extract sleep information, participants were asked to complete a simple sleep diary, as used in the Whitehall study, which consisted of two questions⁴⁶: ‘*what time did you first fall asleep last night?*’ and ‘*what time did you wake up today (eyes open, ready to get up)?*’. Participants were also asked to complete a HETUS time-use diary⁴⁷, and sleep information from here was extracted in cases where data was missing from the simple sleep diary. This multi-instrument groundtruth resulted in 144 activity labels which were then condensed into six free-living behaviour labels {“bicycling”, “mixed”, “sit/stand”, “sleep”, “vehicle”, “walking”} (see mappings at appendix 2a). Fig. S6 shows a visual representation of the structure and time-balance of labels annotated from this free-living dataset. For energy expenditure metabolic equivalent of task (MET) prediction, we used nine behaviour labels {“bicycling”, “sports”, “walking+activity”, “sitstand+activity”, “sitstand+lowactivity”, “sit/stand”, “sleep”, “vehicle”, “walking”} (see mappings at appendix 2b), each with an associated Metabolic Equivalent of Task (MET) score from the compendium of physical activities.

Accelerometer Data Preparation

For data pre-processing we followed procedures used by the UK Biobank accelerometer data processing expert group⁹, that included device calibration⁴⁸, resampling to 100Hz, and removal of noise and gravity^{9,32,33}. For every non-overlapping 30-second time window, we then extracted a 126-dimensional feature vector. Our features are listed in Fig. S2 and were selected from an extensive list of time and frequency domain features described in other studies^{14,17,30,49}. These included: euclidean norm minus one with negative values truncated to zero⁹, it’s mean, standard deviation, coefficient of variation, median, min, max, 25th & 75th percentiles, mean amplitude deviation, mean power deviation, kurtosis & skew, and Fast Fourier Transform (FFT) 1-15Hz. Features also included the following in each axis of movement: mean, range, standard deviation, covariance, and FFT 1-15Hz. Roll, pitch, yaw, x/y/z correlations, frequency and power bands were also extracted.

Activity Classification

For activity classification we use random forests⁵⁰ which offer a powerful nonparametric discriminative method for classification that offers state-of-the-art performance⁵¹. Random forests are collections of decision trees, where each tree is constructed from a training set of feature data with ground truth classes, but each tree observed only a subset of training data and their features. Given the unbalanced nature of our dataset where some behaviours occur rarely, we use balanced Random Forests⁵² to train each tree with a balanced subset of training data. If we have n_{rare} instances of the rarest class, we pick n_{rare} samples, with replacement, of each of our classes to form our training set for each tree. As each tree is given only a small fraction of data, we make many more trees than in a standard random forest so that the same number of data points are sampled in training as with a standard application of random forests. We evaluated different numbers of trees in our random forest, each trained using n_{rare} datapoints from each class of activity. The terminal nodes of each tree contain only one datapoint.

Time Smoothing

Random forests are able to classify datapoints, but do not have an understanding of our data as having come from a time series. Therefore we use a hidden Markov model⁵³ (HMM) to smooth our predictions. The state space of the HMM is the true activity in that time increment and the emissions are the predicted activities from the balanced random forest.

The transition matrix and emission distribution are empirically calculated, the transition matrix from training set sequence of states and the emission probabilities from the out-of-bag class votes of the random forest. Thus the confidence of the random forest in the accuracy of its predictions for an activity follows through into how confident the HMM is for each. With this empirically defined HMM we can then run the Viterbi algorithm⁵⁴ to find the most likely sequence of states given a sequence of observed emissions. This smoothing corrects erroneous predictions from the random forest, such as where the error is a blip of one activity surrounded by another and the transitions between those two classes of activity are rare.

MET Prediction

To predict the MET score we follow the same process of feature extraction, random forest training and HMM definition, but for the nine-class MET-relevant behaviour labels. However, instead of selecting the Viterbi path we obtain the sequence of marginal probabilities for being in each state at each time given the sequence of observations. Each of the nine classes of behaviour is a mix of different activities from the compendium of physical activities, thus a representative MET score was calculated by taking the mean of the MET scores used to construct that class from the training dataset. Finally, the predicted MET score for each 30-second chunk is calculated as the assigned MET scores for each of the 9 states, weighted by the marginal probabilities of being in each of those states.

Extrapolation to large health datasets

We trained a model using all free-living groundtruth data, and applied it to predict behaviour for each 30-second epoch in 103,712 UK Biobank participants' accelerometer data. For any given time window (e.g. one hour, one day, etc.) the probability of a participant engaging in a specific behaviour type was expressed as the number-of-epoch-predictions-for-class divided by the number-of-epochs. Device non-wear time was automatically identified as consecutive stationary episodes lasting for at least 60 minutes²⁶. These non-wear segments of data were imputed with the average of similar time-of-day data points, for each behaviour prediction, from different days of the measurement. We excluded participants whose data could not be calibrated, had too many clipped values⁹, had unrealistically high values (average vector magnitude > 100mg), or who had poor wear-time. We defined minimum wear time criteria as having at least three days (72 hours) of data and also data in each one-hour period of the 24-hour cycle⁹.

Statistical Analysis

To compare machine predicted behaviour from accelerometer data against the free-living groundtruth, we used leave-one-subject-out cross validation, and reported kappa scores for minute-by-minute agreement. The Kappa test reflects the inter-rater agreement between two sources taking into account the likelihood of them agreeing by chance⁵⁵. We used Bland-Altman plots to illustrate daily-level summary agreement between predicted behaviour and the groundtruth. For the UK Biobank dataset, descriptive statistics were used to report accelerometer measured time in {"bicycling", "mixed", "sit/stand", "sleep", "vehicle", "walking"} behaviours. Age groups were categorised into <55, 55-64, and 65+ years. To quantify statistical differences by age, sex, self-rated health, and season, two-way ANOVA linear regression were used, with analysis adjusted for age, sex, ethnicity, area-deprivation, smoking, alcohol, and self-rated-health. Time-of-day (six hour quadrants) and weekday vs. weekend differences in behaviour were reported using two-way repeated measures ANOVA.

As this is used for quantification, rather than hypothesis testing, we report p-values uncorrected for multiple testing. 24-hour activity plots stratified by weekend were used to illustrate accelerometer-measured behavioural profiles. Stack charts were plotted to illustrate the distribution of all objectively measured behavioural types in UK Biobank participants. We used R to perform all statistical analyses⁵⁶.

Data and code availability

Upon publication, the summary variables that we have constructed will be made available as a part of the UK Biobank dataset at <http://biobank.ctsu.ox.ac.uk/crystal/label.cgi?id=1008>. All data processing, feature extraction, machine learning, and analysis code will be available at <https://github.com/activityMonitoring>.

ACKNOWLEDGEMENTS

We would like to thank all participants for agreeing to volunteer in this research. We would like to thank Martin Landray for feedback on this work. For data collection at Oxford and annotation we would like to acknowledge the support of Jonathan Gershuny the UK Economic and Social Research Council (grant number ES/L011662/1), Charlie Foster, Paul Kelly, Teresa Harms, Emma Thomas, Karen Milton, Wong Tsz Yan. The UK Biobank Activity Project and the collection of activity data from participants was funded by the Wellcome Trust (<https://wellcome.ac.uk/>) and the Medical Research Council (<http://www.mrc.ac.uk/>). The analysis was supported by the British Heart Foundation Centre of Research Excellence at Oxford (<http://www.cardioscience.ox.ac.uk/bhf-centre-of-research-excellence>) [grant number RE/13/1/30181 to AD], the Li Ka Shing Foundation (<http://www.lksf.org/>) [to AD]. We would also like to acknowledge the use of the University of Oxford Advanced Research Computing (ARC) facility in carrying out this work. <http://dx.doi.org/10.5281/zenodo.22558> The MRC and Wellcome Trust played a key role in the decision to establish UK Biobank, and the accelerometer data collection. No funding bodies had any role in the analysis, decision to publish, or preparation of the manuscript.

REFERENCES

1. Lee, I.-M. *et al.* Effect of physical inactivity on major non-communicable diseases worldwide: an analysis of burden of disease and life expectancy. *Lancet* **380**, 219–29 (2012).
2. Celis-Morales, C. A. *et al.* Association between active commuting and incident cardiovascular disease, cancer, and mortality: prospective cohort study. *BMJ* **357**, (2017).
3. Shan, Z. *et al.* Sleep Duration and Risk of Type 2 Diabetes: A Meta-analysis of Prospective Studies. *Diabetes Care* **38**, (2015).
4. Kelly, P. *et al.* Systematic review and meta-analysis of reduction in all-cause mortality from walking and cycling and shape of dose response relationship. *Int. J. Behav. Nutr. Phys. Act.* **11**, 132 (2014).
5. Wilmot, E. G. *et al.* Sedentary time in adults and the association with diabetes, cardiovascular disease and death: systematic review and meta-analysis. *Diabetologia* **55**, 2895–2905 (2012).

6. Colbert, L. H., Matthews, C. E., Havighurst, T. C., Kim, K. & Schoeller, D. A. Comparative validity of physical activity measures in older adults. *Med. Sci. Sports Exerc.* **43**, 867–76 (2011).
7. Sudlow, C. *et al.* UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLoS Med.* **12**, e1001779 (2015).
8. Littlejohns, T. J., Sudlow, C., Allen, N. E. & Collins, R. UK Biobank: opportunities for cardiovascular research. *Eur. Heart J.* (2017). doi:10.1093/eurheartj/ehx254
9. Doherty, A. *et al.* Large Scale Population Assessment of Physical Activity Using Wrist Worn Accelerometers: The UK Biobank Study. *PLoS One* **12**, e0169649 (2017).
10. Troiano, R. P., McClain, J. J., Brychta, R. J. & Chen, K. Y. Evolution of accelerometer methods for physical activity research. *Br. J. Sports Med.* **48**, 1019–23 (2014).
11. Menai, M. *et al.* Accelerometer assessed moderate-to-vigorous physical activity and successful ageing: results from the Whitehall II study. *Sci. Rep.* **8**, 45772 (2017).
12. Kerr, J. *et al.* Objective Assessment of Physical Activity: Classifiers for Public Health. *Med. Sci. Sports Exerc.* **48**, 951–7 (2016).
13. Intille, S. S., Lester, J., Sallis, J. F. & Duncan, G. New horizons in sensor development. *Med Sci Sport. Exerc* **44**, 24–31 (2012).
14. Mannini, A., Intille, S. S., Rosenberger, M., Sabatini, A. M. & Haskell, W. Activity recognition using a single accelerometer placed at the wrist or ankle. *Med. Sci. Sports Exerc.* **45**, 2193–203 (2013).
15. Welch, W. A. *et al.* Classification accuracy of the wrist-worn gravity estimator of normal everyday activity accelerometer. *Med. Sci. Sports Exerc.* **45**, 2012–9 (2013).
16. van Hees, V. T., Golubic, R., Ekelund, U. & Brage, S. Impact of study design on development and evaluation of an activity-type classifier. *J. Appl. Physiol.* **114**, 1042–51 (2013).
17. Zhang, S., Rowlands, A. V, Murray, P. & Hurst, T. L. Physical activity classification using the GENE wrist-worn accelerometer. *Med Sci Sport. Exerc* **44**, 742–748 (2012).
18. Staudenmayer, J., Pober, D., Crouter, S., Bassett, D. & Freedson, P. An artificial neural network to estimate physical activity energy expenditure and identify physical activity type from an accelerometer. *J. Appl. Physiol.* **107**, 1300–1307 (2009).
19. Ellis, K., Godbole, S., Kerr, J. & Lanckriet, G. Multi-sensor physical activity recognition in free-living. in *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing Adjunct Publication - UbiComp '14 Adjunct* 431–440 (ACM Press, 2014). doi:10.1145/2638728.2641673
20. Miller, N. E., Welch, W. A., Doherty, A. R. & Strath, S. J. Accuracy of Behavioral Assessment with a Wearable Camera in Semi-Structured and Free Living Conditions in Older Adults. in *American College of Sports Medicine Annual Meeting, 30 May - 03 June* (2017).
21. Montoye, A. H. K., Begum, M., Henning, Z. & Pfeiffer, K. A. Comparison of linear and non-linear models for predicting energy expenditure from raw accelerometer data. *Physiol. Meas.* **38**, 343–357 (2017).
22. Sasaki, J. E. *et al.* Performance of Activity Classification Algorithms in Free-Living Older Adults. *Med. Sci. Sports Exerc.* **48**, 941–50 (2016).
23. Pavey, T. G., Gilson, N. D., Gomersall, S. R., Clark, B. & Trost, S. G. Field evaluation of a random forest activity classifier for wrist-worn accelerometer data. *J. Sci. Med. Sport* (2016). doi:10.1016/j.jsams.2016.06.003
24. Ainsworth, B. E. *et al.* 2011 Compendium of Physical Activities: a second update of codes and MET values. *Med Sci Sport. Exerc* **43**, 1575–1581 (2011).
25. Ellis, K. *et al.* A random forest classifier for the prediction of energy expenditure and type of physical activity from wrist and hip accelerometers. *Physiol. Meas.* **35**, 2191–203 (2014).
26. Doherty, A. *et al.* Large Scale Population Assessment of Physical Activity Using Wrist Worn Accelerometers: The UK Biobank Study. *PLoS One* **12**, e0169649 (2017).

27. Landis, J. R. & Koch, G. C. The measurement of observer agreement for categorical data. *Biometrics* **33**, 159–174 (1977).
28. Doherty, A. R. *et al.* Using wearable cameras to categorise type and context of accelerometer-identified episodes of physical activity. *Int. J. Behav. Nutr. Phys. Act.* **10**, 22 (2013).
29. Kerr, J. *et al.* Using the SenseCam to improve classifications of sedentary behavior in free-living settings. *Am. J. Prev. Med.* **44**, 290–296 (2013).
30. Ellis, K., Kerr, J., Godbole, S., Staudenmayer, J. & Lanckriet, G. Hip and Wrist Accelerometer Algorithms for Free-Living Behavior Classification. *Med. Sci. Sports Exerc.* **48**, 933–40 (2016).
31. Bonomi, A. G., Plasqui, G., Goris, A. H. C. & Westerterp, K. R. Improving assessment of daily energy expenditure by identifying types of physical activity with a single accelerometer. *J. Appl. Physiol.* **107**, (2009).
32. Sabia, S. *et al.* Association Between Questionnaire- and Accelerometer-Assessed Physical Activity: The Role of Sociodemographic Factors. *Am. J. Epidemiol.* **179**, 781–790 (2014).
33. da Silva, I. C. *et al.* Physical activity levels in three Brazilian birth cohorts as assessed with raw triaxial wrist accelerometry. *Int. J. Epidemiol.* **43**, 1959–68 (2014).
34. Borazio, M., Berlin, E., Kucukyildiz, N., Scholl, P. & Laerhoven, K. Van. Towards Benchmarked Sleep Detection with Wrist-Worn Sensing Units. in *2014 IEEE International Conference on Healthcare Informatics* 125–134 (IEEE, 2014). doi:10.1109/ICHI.2014.24
35. Strath, S. J. *et al.* Guide to the assessment of physical activity: Clinical and research applications: a scientific statement from the American Heart Association. *Circulation* **128**, 2259–79 (2013).
36. White, T., Westgate, K., Wareham, N. J. & Brage, S. Estimation of physical activity energy expenditure during free-living from wrist accelerometry in UK adults. *PLoS One* (in press), (2016).
37. Kelly, P. *et al.* Developing a Method to Test the Validity of 24 Hour Time Use Diaries Using Wearable Cameras: A Feasibility Pilot. *PLoS One* **10**, e0142198 (2015).
38. Ladha, C., Ladha, K., Jackson, D. & Olivier, P. Shaker Table Validation Of Openmovement Ax3 Accelerometer. in *3rd International Conference on Ambulatory Monitoring of Physical Activity and Movement* 69–70 (2013).
39. Eslinger, D. W. *et al.* Validation of the GENEA Accelerometer. *Med Sci Sport. Exerc* **43**, 1085–1093 (2011).
40. Doherty, A. R. *et al.* Wearable cameras in health: The state of the art and future possibilities. *Am. J. Prev. Med.* **44**, 320–323 (2013).
41. Hodges, S. *et al.* SenseCam: A Retrospective Memory Aid. in *UbiComp: 8th International Conference on Ubiquitous Computing* **4602**, 177–193 (Springer, 2006).
42. Kelly, P. *et al.* High group level validity but high random error of a self-report travel diary, as assessed by wearable cameras. *J. Transp. Heal.* (2014). doi:10.1016/j.jth.2014.04.003
43. Doherty, A. R. *et al.* Using wearable cameras to categorise type and context of accelerometer-identified episodes of physical activity. *Int. J. Behav. Nutr. Phys. Act.* **10**, (2013).
44. Kelly, P. *et al.* Ethics of using wearable cameras devices in health behaviour research. *Am J Prev Med* **44**, 314–319 (2013).
45. Doherty, A. R., Moulin, C. J. A. & Smeaton, A. F. Automatically assisting human memory: A SenseCam browser. *Memory* **19**, 785–795 (2011).
46. van Hees, V. T. *et al.* A Novel, Open Access Method to Assess Sleep Duration Using a Wrist-Worn Accelerometer. *PLoS One* **10**, e0142533 (2015).
47. Eurostat. *Harmonised European Time Use Surveys: 2008 Guidelines.* (2008).
48. van Hees, V. T. *et al.* Autocalibration of accelerometer data for free-living physical activity assessment using local gravity and temperature: an evaluation on four continents. *J. Appl. Physiol.* **117**, 738–44 (2014).

49. Vähä-Ypyä, H., Vasankari, T., Husu, P., Suni, J. & Sievänen, H. A universal, accurate intensity-based classification of different physical activities using raw data of accelerometer. *Clin. Physiol. Funct. Imaging* **35**, 64–70 (2015).
50. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
51. Fernández-Delgado, M., Cernadas, E., Barro, S. & Amorim, D. Do we Need Hundreds of Classifiers to Solve Real World Classification Problems? *J. Mach. Learn. Res.* **15**, 3133–3181 (2014).
52. Chen, C., Liaw, A. & Breiman, L. *Using random forest to learn imbalanced data*. University of California, Berkeley (2004).
53. Rabiner, L. & Juang, B. An introduction to hidden Markov models. *IEEE ASSP Mag.* **3**, 4–16 (1986).
54. Forney, G. D. The viterbi algorithm. *Proc. IEEE* **61**, 268–278 (1973).
55. Viera, A. J. & Garrett, J. M. Understanding interobserver agreement: the kappa statistic. *Fam. Med.* **37**, 360–3 (2005).
56. R Core Team. R: A Language and Environment for Statistical Computing. (2016).

Tables

Table 1. Percentage of machine-learned behaviours automatically classified from wrist-worn accelerometer data. Confusion matrix after leave-one-out validation on 84,616 labelled minutes of human activity in free-living environments: the CAPTURE-24 study 2014-2015 (n = 57).

Prediction→ Ground truth↓	Sleep	Sit/stand	Vehicle	Walking	Mixed- activity	Bicycling
Sleep	95%	4%	<1%	<1%	1%	<1%
Sit/stand	4%	85%	1%	2%	9%	0%
Vehicle	<1%	8%	81%	4%	7%	0%
Walking	<1%	13%	1%	51%	33%	2%
Mixed-activity	<1%	11%	3%	11%	75%	1%
Bicycling	<1%	<1%	3%	15%	7%	75%

Table 2. Objective machine-learned measures of physical activity (vector magnitude), sleep, walking, sitting-or-standing, bicycling, vehicle, and mixed activity time: the UK Biobank study 2013-2015 (n = 96,609).

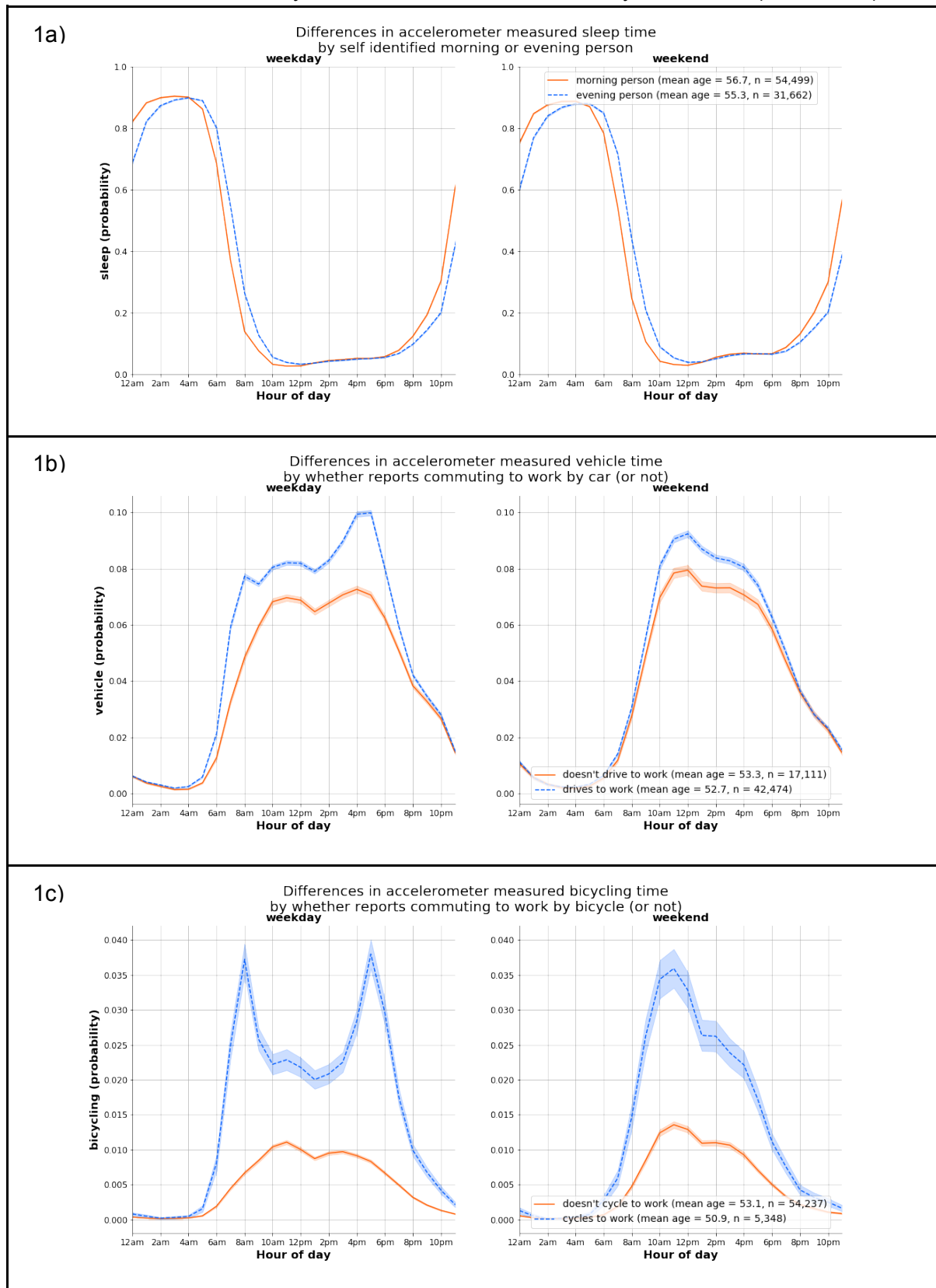
	Individuals	Vector magnitude	MET	sleep	walking	sit/stand	bicycling	vehicle	mixed
	[n]	[mg]	[MET hrs/day]	[% time]					
				[mean ± stdev]					
Age									
<55	20,543	31.2 ± 9.0	39.5 ± 3.0	33.6 ± 6.0	7.3 ± 3.5	33.1 ± 7.9	0.6 ± 1.1	4.6 ± 2.8	15.2 ± 6.5
55-64	33,874	28.9 ± 8.2	39.4 ± 2.8	34.4 ± 5.9	6.9 ± 3.6	33.3 ± 7.8	0.5 ± 1.1	4.5 ± 2.8	15.5 ± 6.3
65+	42,192	25.9 ± 7.3	39.0 ± 2.8	35.2 ± 5.9	6.1 ± 3.4	34.3 ± 7.7	0.5 ± 0.9	4.1 ± 2.5	15.8 ± 6.2
p value ^A		<1x10 ⁻³⁰⁰	2x10 ⁻¹²³	3x10 ⁻²³³	<1x10 ⁻³⁰⁰	1x10 ⁻¹⁰³	2x10 ⁻⁷⁰	2x10 ⁻¹²⁴	2x10 ⁻²⁶
Cohen's d		.65	.18	.27	.34	.15	.14	.17	.09
Sex									
Women	54,396	28.5 ± 8.0	39.7 ± 2.8	34.6 ± 5.7	6.2 ± 3.2	32.5 ± 7.4	0.3 ± 0.8	4.1 ± 2.5	17.4 ± 6.2
Men	42,213	27.5 ± 8.6	38.7 ± 2.9	34.5 ± 6.3	7.2 ± 3.8	35.3 ± 8.0	0.7 ± 1.3	4.6 ± 2.9	13.2 ± 5.7
p value ^A		2x10 ⁻⁴³	<1x10 ⁻³⁰⁰	4x10 ⁻⁴	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁸	<1x10 ⁻³⁰⁰	3x10 ⁻¹⁷¹	<1x10 ⁻³⁰⁰
Cohen's d		.12	.33	.01	.29	.35	.36	.17	.70
Self-rated health									
Excellent	21,201	30.4 ± 8.8	39.6 ± 2.8	34.3 ± 5.7	7.2 ± 3.5	33.0 ± 7.4	0.6 ± 1.2	4.5 ± 2.6	15.8 ± 6.2
Good	58,024	28.2 ± 7.9	39.3 ± 2.8	34.6 ± 5.9	6.7 ± 3.5	33.5 ± 7.6	0.5 ± 1.0	4.3 ± 2.6	15.7 ± 6.3
Fair/poor	17,848	25.3 ± 7.8	38.6 ± 3.0	34.8 ± 6.7	5.9 ± 3.6	35.3 ± 8.6	0.4 ± 0.9	4.1 ± 2.8	14.7 ± 6.5
p value ^A		<1x10 ⁻³⁰⁰	4x10 ⁻²²⁹	5x10 ⁻¹⁸	6x10 ⁻²⁹⁹	2x10 ⁻¹⁴¹	1x10 ⁻¹³⁰	8x10 ⁻³²	1x10 ⁻⁴⁹
Cohen's d		.61	.36	.09	.37	.29	.24	.14	.18
Time of day									
0-5.59 _{am}	96,609	4.7 ± 3.7	24.5 ± 2.3	85.8 ± 13.9	0.4 ± 1.1	5.6 ± 7.4	0.0 ± 0.3	0.4 ± 1.3	1.4 ± 2.3
6-11.59 _{am}	96,609	38.2 ± 15.3	43.6 ± 5.3	27.5 ± 12.9	9.2 ± 6.1	30.7 ± 11.0	0.8 ± 1.9	5.3 ± 4.0	21.8 ± 9.8
12-5.59 _{pm}	96,609	43.8 ± 14.8	48.6 ± 4.9	4.6 ± 5.9	12.1 ± 6.9	46.8 ± 13.2	1.0 ± 2.1	8.0 ± 5.2	23.8 ± 11.6
6-11.59 _{pm}	96,609	25.8 ± 10.8	40.3 ± 4.3	20.2 ± 12.3	4.8 ± 3.9	51.7 ± 13.3	0.3 ± 0.9	3.7 ± 3.7	15.3 ± 8.0
p value ^B		<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰
Cohen's d		3.6	6.3	7.6	2.4	4.3	.66	2.0	2.7
Day									
Weekday	96,609	28.2 ± 8.5	39.4 ± 3.0	34.1 ± 6.1	6.8 ± 3.7	34.3 ± 8.3	0.5 ± 1.1	4.5 ± 3.0	15.6 ± 6.7
Weekend	96,609	27.7 ± 9.8	38.8 ± 3.8	35.6 ± 8.6	6.3 ± 4.4	32.3 ± 10.3	0.5 ± 1.3	3.9 ± 3.4	15.5 ± 7.5
p value ^B		4x10 ⁻⁹⁸	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	<1x10 ⁻³⁰⁰	4x10 ⁻¹⁹	<1x10 ⁻³⁰⁰	5x10 ⁻¹¹
Cohen's d		.05	.16	.19	.14	.22	.03	.17	.02
Season									
Spring	21,925	28.6 ± 8.4	39.4 ± 2.8	34.4 ± 5.9	6.9 ± 3.6	33.5 ± 7.8	0.6 ± 1.1	4.3 ± 2.7	15.5 ± 6.3
Summer	25,382	28.6 ± 8.4	39.6 ± 2.9	34.1 ± 5.9	6.6 ± 3.5	33.6 ± 7.9	0.6 ± 1.2	4.3 ± 2.7	16.1 ± 6.5
Autumn	28,819	27.9 ± 8.2	39.2 ± 2.8	34.6 ± 6.0	6.6 ± 3.5	33.9 ± 7.7	0.5 ± 1.0	4.4 ± 2.7	15.4 ± 6.3
Winter	20,483	27.2 ± 7.9	38.9 ± 2.8	35.2 ± 6.1	6.4 ± 3.4	33.8 ± 7.8	0.4 ± 0.9	4.2 ± 2.6	15.2 ± 6.2
p value ^A		1x10 ⁻¹⁰²	1x10 ⁻¹⁵⁰	2x10 ⁻⁸³	4x10 ⁻⁵¹	6x10 ⁻⁰⁶	6x10 ⁻¹⁰⁶	3x10 ⁻¹⁴	1x10 ⁻⁵⁷
Cohen's d		.18	.23	.18	.15	.05	.19	.07	.13

^A Age, sex, self-rated health, season (Spring starting on 1st March): Two-way analysis of variance test used to compare metrics between groups adjusting for age, sex, ethnicity, area-deprivation, smoking, alcohol, and self-rated health.

^B Time of day, day: Repeated two-way analysis of variance test used to compare metrics within individuals and between groups adjusting for age, sex, ethnicity, area-deprivation, smoking, alcohol, and self-rated health.

Figures

Fig. 1. Variation in accelerometer-measured behaviour types across the day by participant characteristics and weekday/weekend: the UK Biobank study 2013-2015 (n = 96,609).



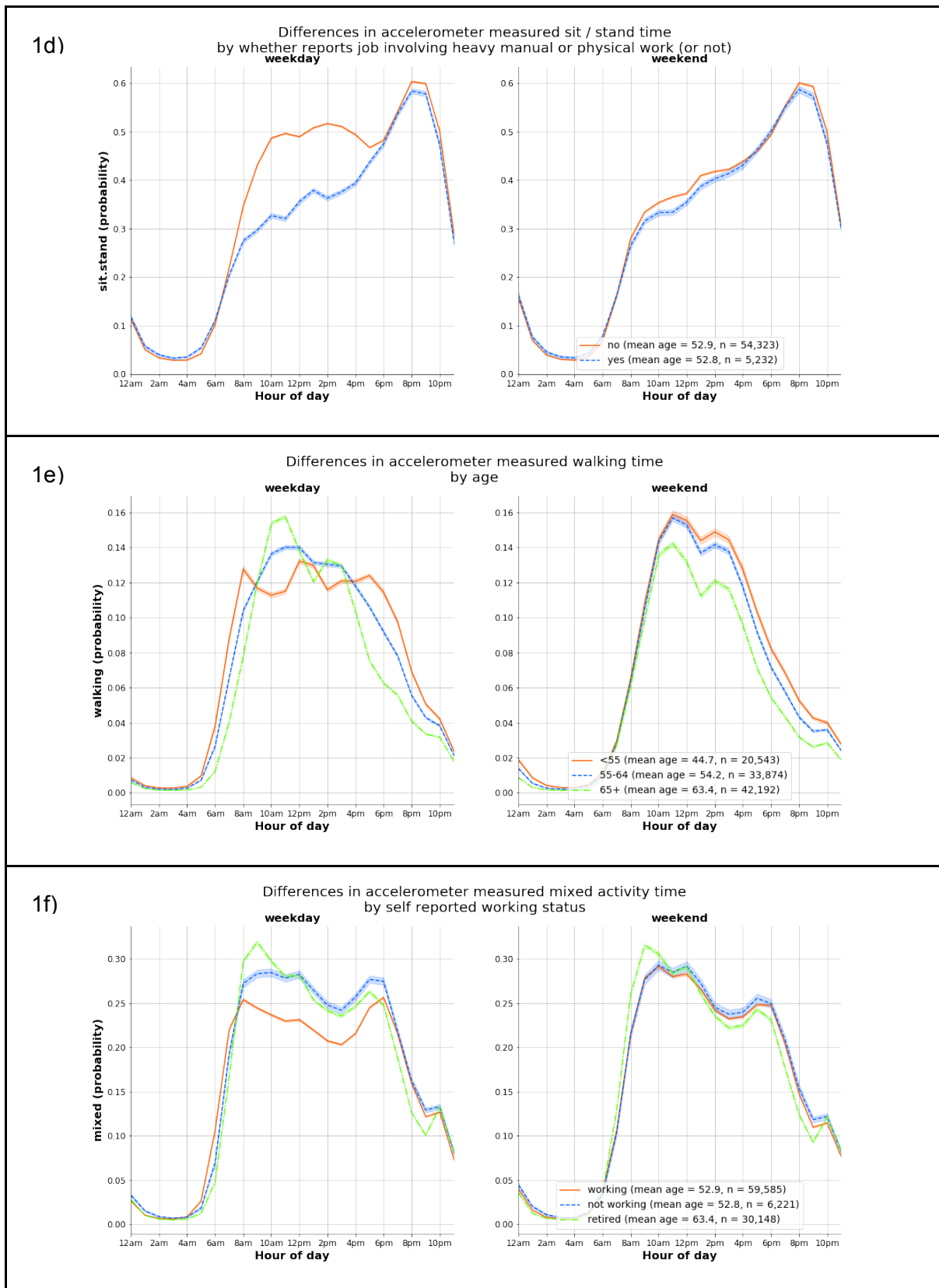


Fig. 2. Variation in accelerometer-measured time by activity type: the UK Biobank study 2013-2015 (n = 96,609).

