

Smart-Plexer: a breakthrough workflow for hybrid development of multiplex PCR assays

Luca Miglietta

Imperial College London <https://orcid.org/0000-0002-6652-4785>

Yuwen Chen

Imperial College London

Zhi Luo

Imperial College London <https://orcid.org/0000-0002-7567-1082>

Ke Xu

Imperial College London

Ning Ding

Imperial College London

Tianyi Peng

Imperial College London

Ahmad Moniri

Imperial College London

Louis Kreitmann

Imperial College London

Miguel Cacho-Soblechero

Imperial College London

Alison Holmes

Imperial College London

Pantelis Georgiou

Imperial College London <https://orcid.org/0000-0003-2476-3857>

Jesus Rodriguez-Manzano (✉ j.rodriguez-manzano@imperial.ac.uk)

Imperial College London

Article

Keywords:

Posted Date: June 17th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1765213/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Smart-Plexer: a breakthrough workflow for hybrid development of multiplex PCR assays

Luca Miglietta^{1,2}, Yuwen Chen^{2,5}, Zhi Luo^{1,5}, Ke Xu^{1,2}, Ning Ding¹, Tianyi Peng², Ahmad Moniri², Louis Kreitmann¹, Miguel Cacho-Soblechero², Alison Holmes¹, Pantelis Georgiou² and Jesus Rodriguez-Manzano^{1*}

¹Department of Infectious Disease, Faculty of Medicine, Imperial College London, London, UK.

²Department of Electrical and Electronic Engineering, Faculty of Engineering, Imperial College London, London, UK.

⁵These authors contributed equally.

*Corresponding author: j.rodriquez-manzano@imperial.ac.uk

Developing multiplex PCR assays requires an extensive amount of experimental testing, the number of which exponentially increases by the number of multiplexed targets. Dedicated efforts must be devoted to the design of optimal multiplex assays for specific and sensitive identification of multiple analytes in a single well reaction. Inspired by data-driven approaches, we reinvent the way of designing and developing multiplex assays by proposing a hybrid, easy-to-use workflow, named Smart-Plexer, which couples empirical testing of singleplex assays and computer simulation of multiplexing. The Smart-Plexer leverages kinetic inter-target distances among amplification curves to generate optimal multiplex PCR primer sets for accurate multi-pathogen identification. The optimal single-channel assays, together with a novel data-driven approach, Amplification Curve Analysis (ACA), were demonstrated to be capable of classifying the presence of desired targets in a single test for seven common respiratory infection pathogens.

Quantitative Polymerase Chain Reaction (qPCR) allows to continuously monitor the kinetic signature of a specific amplification event due to the mutual interaction of oligonucleotides and their specific template¹⁻³. The extraordinary ease and reliability of this golden standard method for Nucleic Acid Amplification Tests (NAATs) have improved routine diagnostics in several fields and, more recently, played a crucial role during the COVID-19 pandemic, one of the ten deadliest infectious diseases in history⁴⁻⁶. This epidemic has further highlighted the need for more cost-effective and provisional diagnoses, and for enhancing the diagnostic capabilities of conventional instruments along with point-of-care devices⁷⁻¹⁰. As the pandemic comes to an end, the focus on developing NAATs for simultaneous detection of multiple respiratory pathogens alongside COVID-19 has been drastically increased¹¹⁻¹³. There is an emerging demand for rapid, affordable, and reliable molecular tests for multiple identifications of infectious disease^{7,9}.

Current screening strategies of multiple pathogens are reported to be expensive, sample consuming and, in some cases, inaccurate¹⁴⁻¹⁶. As a result, multiplex PCR is emerging as an inexpensive alternative for multi-target identification¹⁷⁻²⁰. Many efforts have been made in developing novel methods to increase the number of targets detected by multiplex assays and to enhance the accurate identification of multiple infectious sources in a single test^{21,22}. Advances in multi-pathogen detection include the use of High-Resolution Melting Analysis (HRMA), fluorescent probe-based method, or restriction enzyme digestion²³⁻²⁵. Recently, the emergence of machine learning approaches in clinical diagnostics has highlighted the potential of data-driven multiplexing, which, compared to conventional methods, unbars limitations in terms of throughput, costs, time and reliability^{26,27}. A few methods have been proposed using either melting curve analysis (intercalating dye-based chemistries) or the final fluorescence intensity (probe-based assays) as features for machine learning algorithms^{28,29}. Moreover, using cutting-edge signal processing and tailored amplification chemistries, state-of-the-art identification performance has been achieved by leveraging the kinetic information encoded in the entire amplification curve from multiplex PCR assays. A novel learning-based methodology called Amplification Curve Analysis (ACA) has been recently reported as a digital tool to expand multiplex capabilities of real-time PCR-based diagnostic platforms, increasing the number of detectable targets per fluorescent channel in a single reaction without hardware modification³⁰⁻³³.

However, the development of multiplex PCR assays is still restrained as extensive experimental testing are required to assess the assay's analytical performance, such as cross-reactivity, specificity, and sensitivity^{21,34-36}. One of the biggest challenges in multiplexing is the complexity of assay design, which dramatically increases with the number of targets, making the development costly, lengthy and resource consuming in

53 the wet laboratory^{14,37}. For N_t multiplexed targets, if N_{ps} candidate primer sets are designed for each of
54 them (which is trivial progress for well-designed singleplex assays), the total number of possible multiplex
55 assay combinations is $N_c = N_{ps}^{N_t}$ (e.g. $N_c = 16,384$ when $N_{ps} = 4$ and $N_t = 7$). The N_c increases
56 exponentially with N_t , making it impractical to find the optimal combination by wet-lab experiments in high-
57 level multiplexing. Therefore, an *in-silico* simulation method is required for fast screening and for narrowing
58 down selections of multiplex assays.

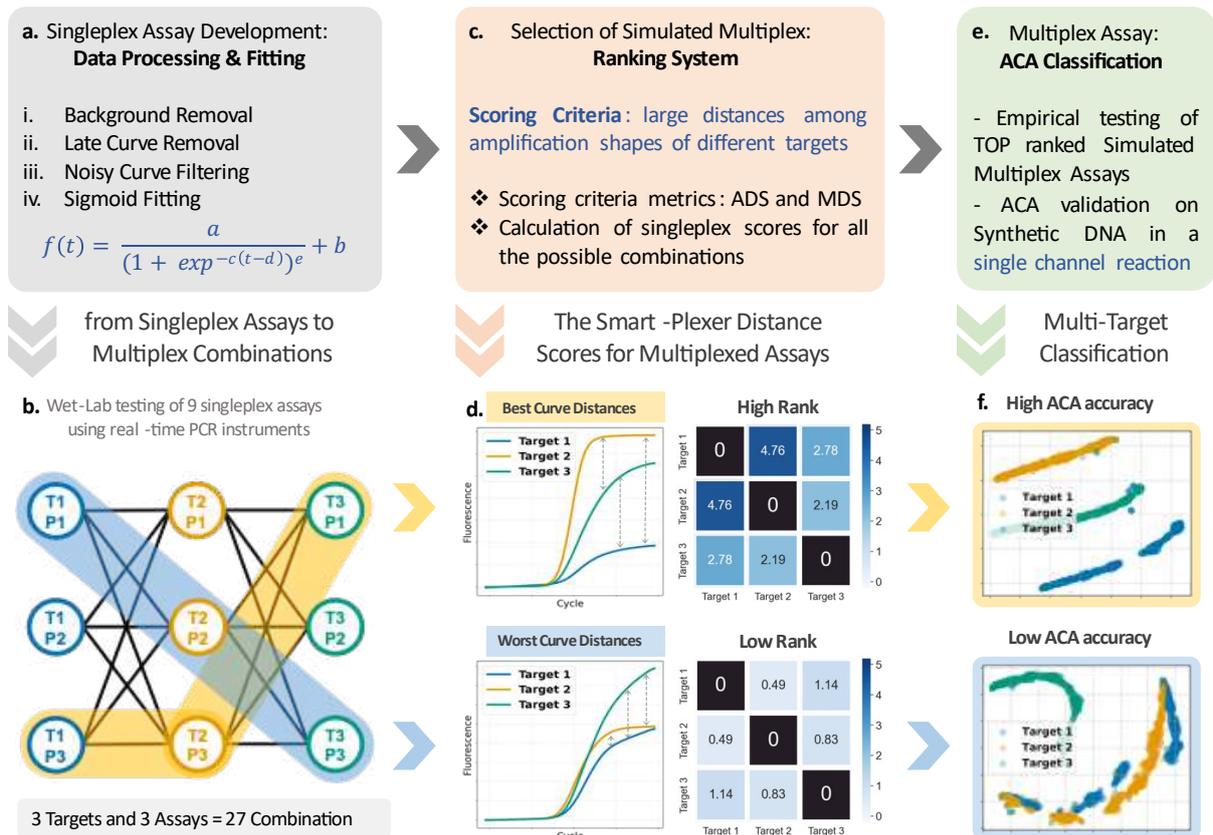
59 To address this problem, here we present the Smart-Plexer, a mathematical algorithm capable of
60 simulating thousands of possible multiplex assay combinations based on singleplex real-time digital PCR
61 (qPCR) data. We aim to demonstrate the use of this new methodology by developing a TaqMan-based
62 multiplex assay, in a single fluorescent channel, for the specific and sensitive detection of seven common
63 respiratory tract infection (RTI) pathogens. This work is two-folded: First, we validated the Smart-Plexer by
64 comparing the performance of all possible simulated and empirical combinations in 3-plex, showing a strong
65 correlation between *in-silico* and lab-tested multiplexes; Second, we assessed the proposed pipeline in high-
66 level multiplex (7-plex) by evaluating the ACA classification performance on synthetic DNA and clinical
67 samples. We demonstrated that, out of 4,608 simulated combinations, an optimal multiplex assay could be
68 developed using this novel framework to detect seven common respiratory pathogens accurately in qPCR.

69 Results

70 **Smart-Plexer design framework.** We developed the Smart-Plexer, a framework that uses singleplex PCR
71 reactions as a 'card deck' to generate a 'winning combination' of the multiplex assay. After deciding the
72 number of targets intended to multiplex, the Smart-Plexer takes as input a dataset generated from real-time
73 PCR reactions with a single primer set (or singleplex assay) and a single target. Given the desired number of
74 targets to multiplex in a single channel PCR, sigmoidal curves generated from all the singleplex/target
75 interactions can be combined to simulate curves from a multiplex assay (**Fig. 1**). These simulations of assay
76 combinations are then empirically tested in wet-lab multiplex tests for each target to evaluate changes in
77 the curve shape of the amplification reaction during the transition from singleplex to multiplex environment
78 (empirical multiplex). Moreover, to identify multiple targets with empirical multiplexes, this framework was
79 coupled and evaluated with the ACA methodology.

80 As the ACA is a classifier recognising clusters from different amplification shapes (which, in our case,
81 represent different targets), it is crucial to maintain differences among sigmoidal trends *in silico*. Therefore,
82 those differences across targets can be computed using the Smart-Plexer method through distance
83 measurements (such as Euclidian distance). This novel framework is capable of distance calculation from
84 either the entire amplification curve or its sigmoidal features. The average of computed distances among all
85 the targets is used to rank each combination of singleplex (or simulated multiplex) from high to low inter-
86 curve similarity values. Moreover, the ranking system takes the minimum distance between the two closest
87 targets to ensure that simulated multiplex with high average values is not dependent on the high difference
88 of only a group of curves. When two amplification curves have high similarity, hence a small distance value,
89 the ACA classifier will not work efficiently to identify either target. Therefore, the rank of the combination
90 depends on both average and minimum distance scores. A set of singleplex assays from the top ranks were
91 selected as simulated multiplex for the empirical validation in the laboratory, and the ACA performance was
92 assessed.

93 To compute distances between amplification curves, the Smart-Plexer requires a filtering process where
94 the amplification data generated undergo the following steps: (i) subtraction of curve background to remove
95 the fluorescence signal noise at the starting cycles, (ii) removal of late amplification curves to exclude non-
96 plateau reactions, (iii) removal of noisy curves to exclude non-sigmoidal shapes resulted by operator or
97 instrumentation faults³⁸. The following step comprised of a fitting equation using the 5-parameter model
98 proposed by Spiess et al.³⁹



99
100
101
102
103
104
105
106
107
108
109
110

Fig. 1 | Smart-Plexer workflow. **a.** Given a dataset of singleplex real-time PCR reactions (real-time amplification curves), a processing step is applied (a.i-a.iii). The processed curves are fitted following the equation depicted in step a.iv. An example is given in b., where each curve resulting from singleplex reactions is used in a simulation of multiplex assays. Three targets are considered, and each of them has three unique singleplex assays (a total of 27 simulated combinations). **c.** The simulated multiplex scores are calculated from the Smart-Plexer according to the Scoring Criteria. **d.** Distances within curves from different targets are calculated based on mathematical algorithms (such as Euclidean), and as shown in the confusion matrices, resulting values are used to rank multiplex assays from high (high distances within targets) to low (low distances within targets). **e.** High-rank multiplex assays are chosen for empirical testing, and the ACA method is used to evaluate the classification performance on target identification of each selected multiplex. **f.** Cluster visualisation with 2-D t-SNE represents the difference in inter-target distances between a High-Rank and a Low-Rank multiplex, resulting in high and low ACA classification accuracy, respectively.

111
112
113
114
115
116
117
118

Selection of representative amplification curve. The ACA method uses the entire amplification curve as a time series where fluorescence values change as the number of cycles increases. Firstly, we chose the entire raw amplification curve generated from the real-time PCR reaction as the input of the Smart-Plexer. Secondly, the framework was evaluated using curves normalised with the final fluorescence intensity (FFI) as input to assess performance changes by removing the absolute fluorescence information. To further investigate changes related to different curve representations and different levels of data abstractions (feature dimensions) provided to the Smart-Plexer, sigmoidal parameters generated from a fitting model were also used as input to assess the influence on this framework.

119
120
121
122
123
124
125
126
127

To evaluate the best fitting model, primary efforts have been focused on the selection of an appropriate equation. Several methods have been proposed to efficiently model the real-time PCR sigmoid, such as four, five, and six-parametric functions^{40,39,41}. As a case study, we retrieved the amplification curve data previously reported by Moniri et al., 2020³⁰. Using raw curves as input, after sigmoidal fitting, we calculated the Mean Square Error (MSE) between the raw and the fitted curves for the entire dataset. As shown in **Supplementary Table 1**, the lowest MSE is achieved with the five-parametric model (MSE=0.0036). The rising MSE in six-parameter sigmoid fitting is caused by unsuccessful optimisation resulting from a larger searching dimension. Based on the lowest MSE value, it is determined to utilise the five-parameter sigmoid function to extract features, and the equation is given below:

128

$$f(t) = \frac{a}{(1 + \exp^{-c(t-d)})^e} + b$$

129 where t is the amplification time (or PCR cycle), $f(t)$ is the fluorescence at time t , a is the maximum
130 fluorescence, b is the baseline of the sigmoid, c is related to the slope of the curve, d is the fractional cycle
131 of the inflection point, and e allows for an asymmetric shape (Richard's coefficient).

132 The three different curve representations (raw curves, FFI normalised curves and fitted parameters) were
133 further used to evaluate the transferability from singleplex to multiplex reactions in the Smart-Plexer.
134

135 **Average Distance Score (ADS) and Minimum Distance Scores (MDS) based on curve distances to rank**
136 **multiplex assays.** We developed two distance metrics to measure transferability from simulated to empirical
137 multiplexes, since it is hypothesised that distances between amplification curves should be maintained
138 during the transition from singleplex to multiplex environments.

139 It is possible to calculate distances between two distinct curves by considering them as two data points in
140 the multidimensional space and quantify their distances using various metrics (i.e., Euclidian, Cosine and
141 Manhattan). In a single channel multiplex assay, the number of primer sets present in the reaction equals
142 the number of targets (N_t), therefore the number of distances (N_d) among curves of different targets is
143 represented by the following formula:

$$144 \quad N_d = \binom{N_t}{2} = \frac{N_t(N_t - 1)}{2}$$

145 The average of all the distances is used to assign a score to the multiplex assay called Average Distance
146 Score (ADS). The ADS provides information on the overall distances across targets, and the higher its values
147 are, the more distant the curves are, and better ACA performance is expected (as distances are related to
148 data point clusters). A high ADS does not guarantee a large distance between every two targets of the
149 multiplex. To overcome this limitation, we considered a second metric called Minimum Distance Score
150 (MDS), the distance value of the two closest curves (minimum value of the given N_d distances).

151 The ADS and MDS narrow down the selection of empirical testing for the highest performing multiplexes
152 using a ranking system. Moreover, they are used to validate that inter-curve distance information is
153 maintained during the transition from simulated to empirical multiplexes, and they can be used to develop
154 assays *in silico* more suitable for ACA, skipping costly and timely laboratory testing.
155

156 **Smart-Plexer validation using a 3-plex assay.** To assess the performance of the Smart-Plexer for both *in-silico*
157 multiplex development and ACA classification accuracy, we designed three primer sets for three selected
158 targets using synthetic DNA and tested them in real-time digital PCR (qdPCR): Adenovirus (HAdV), Human
159 coronavirus HKU1 (HCoV-HKU1) and Middle East respiratory syndrome-related coronavirus (MERS-CoV). As
160 shown in **Fig. 1**, the number of combinations to test using N_t targets ($N_t = 3$) and N_{p_s} assays for each target
161 ($N_{p_s} = 3$) is 27 ($N_c = N_{p_s}^{N_t} = 27$ combinations, listed in **Supplementary Table 2**). Three targets were
162 chosen to validate the Smart-Plexer because a complete comparison of all the 27 simulated and empirical
163 multiplex assays can be experimentally conducted as the number of wet-lab experiments is achievable
164 ($N_c \times N_t = 81$ tests).

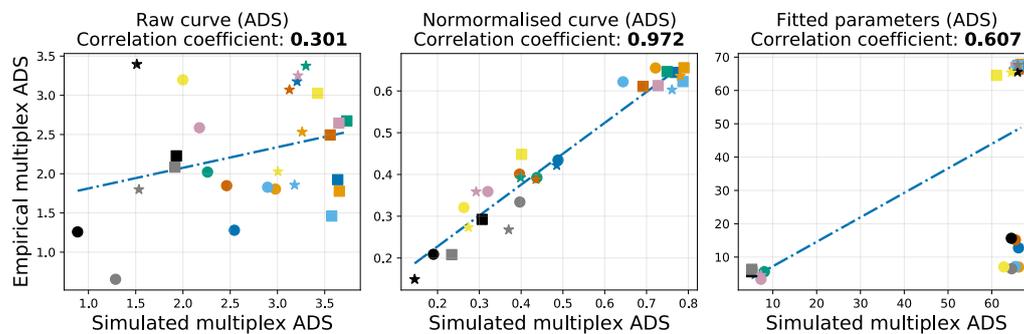
165 The wet-lab testing of each primer set (or singleplex assay) was conducted, and the resulting raw data were
166 combined in a total of 27 simulated multiplexes as explained before. Similarly, experiments were carried out
167 on combinations of primer sets (or empirical multiplex assays) in a single channel reaction. A group of
168 amplification curves, which can be considered as data points in multidimensional spaces, were generated
169 from a unique interaction between each assay and its specific target. The median of these data points was
170 calculated to represent each group of curves. Furthermore, distances among all the curve medians were
171 used to generate the ADS and MDS of all the possible combinations **Fig. 2a-b** visually represent the
172 correlation between the *in-silico* and wet-lab tested assays using ADS and MDS in simulated and empirical
173 multiplexes. Pearson coefficients were reported for both ADS as 0.301, 0.972 and 0.607, and MDS as 0.092,
174 0.761 and 0.686, for raw curve, normalised curve and fitted parameters, respectively (visual representations
175 of each curve type/parameters are depicted in **Fig. 2c**, and ADS and MDS for all the curve types/combinations
176 are reported in **Supplementary Table 3**).

177 It can be observed that normalised curve correlations scored higher than the rest in both ADS and MDS,
178 showing that simulated and empirical multiplex are correlated if FFI is discarded. It is also important to note
179 that the use of all the five curve parameters worsens the correlation as the bimodal distribution of parameter

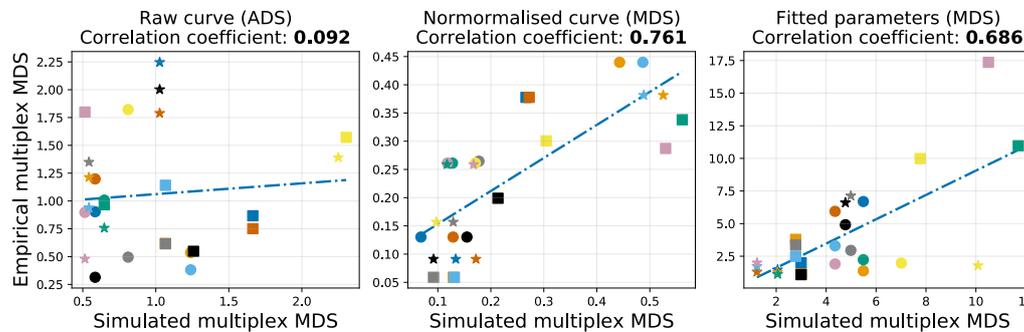
180 “e” negatively influences the correlation, as discussed by Miglietta et al., 2022³⁸. Moreover, the correlation
 181 from singleplex to multiplex might be affected by the fact that the “d” parameter is related to the cycle
 182 threshold (C_t) of the amplification curve. Target concentration can be influenced by instrumentation,
 183 operator, and experimental errors; therefore, variabilities of C_t can easily mislead the correlation of the five
 184 parameters using “d”. Moreover, the scope of conducting this correlation is to compare purely sigmoidal
 185 shapes, and concentrations of the nucleic acid targets should not affect the distance values of two curves.
 186 In addition, the use of parameter “a” and “b” is redundant as: (i) “a” is related to the FFI, and as shown in
 187 the middle plot of Fig. 2a-b, FFI is not relevant to the distance correlation and (ii) all curves present in this
 188 dataset were processed with a background removal (baseline correction) and all “b” parameters were
 189 levelled to almost zero.

190 These discoveries on the correlation between simulated and empirical multiplex distances inspired us to
 191 seek a more representative feature which would maintain the information of distances during the
 192 translation from a singleplex to a multiplex environment. As mentioned before, the parameter “a”, “b”, “d”
 193 and “e” can negatively influence the correlation for both ADS and MDS; therefore, we focus on the “c”
 194 parameter.

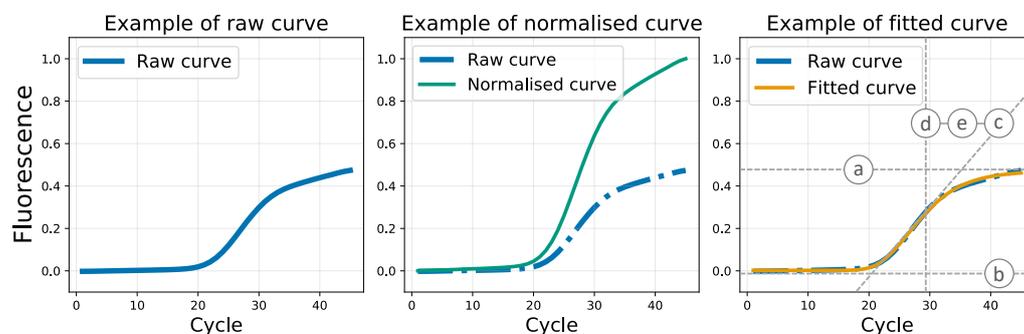
a) ADS of 3-plex experiment



b) MDS of 3-plex experiment



c) Curve type of 3-plex experiment



195
 196 **Fig. 2 | Representative features investigation based on the 3-plex assay.** a. The correlations of the Average distance score
 197 (ADS) between simulated and empirical multiplexes for the three types of curves/parameters (Raw curve, normalised curve
 198 and fitted parameters) are presented (from left to right in the same order). For each plot, each point with unique colour and
 199 shape corresponds to combination 1 to 27. The blue dashed lines are computed using linear regression. The Pearson

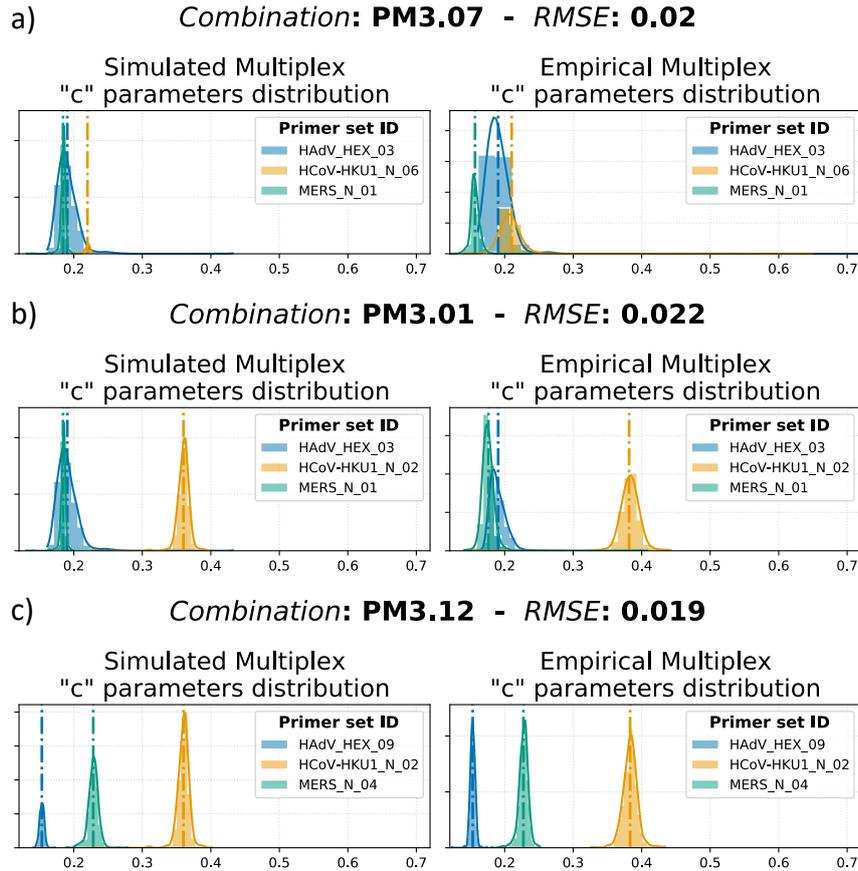
200 coefficients for all three plots are calculated. **b.** Similarly, the correlations of Minimum distance score (MDS) are depicted for
 201 the three curve representations. **c.** Illustration of the three types of curve representations. Examples of raw amplification curve
 202 (after data processing), normalised curve (computed based on the FFI) and fitted curve/parameters are presented from left
 203 to right. The fitted curve is computed with a 5-parameter Sigmoid function using raw curves. As a result of this, we can obtain
 204 both fitted parameters (“a”, “b”, “c”, “d”, “e”) and fitted curve (predicted fluorescence values corresponding to each cycle
 205 from the 5-parameter Sigmoid model with fitted parameters).
 206

207 **The key parameter for curve distances correlation in multiplex assays: the “slope”.** The previous section
 208 reported all the correlation coefficients for ADS and MDS between simulated and empirical multiplexes, in
 209 concomitance with different curve representations: raw curves, normalised curves, and fitting parameters.
 210 Both ADS and MDS showed the maximum correlation values when considering normalised curves. Those
 211 results, along with our discussion on the fitted parameters in the previous section, indicate that reducing
 212 the information contained in the amplification curve is beneficial. This section explores how the “c”
 213 parameter preserves distance information from singleplex to multiple environments of each primer
 214 set/target reaction.

215 In the 3-plex validation, each singleplex assay was tested against its specific target (N=9), resulting in 27
 216 different combinations of simulated multiplexes. Moreover, the “c” parameters were fitted and extracted
 217 from 27 empirically tested multiplex assays (81 tests). **Supplementary Fig. 1** shows the correlation between
 218 simulated and empirical ADS and MDS calculated from “c” parameters with correlation coefficients of 0.973
 219 and 0.774, respectively. To further evaluate whether “c” parameter distributions were maintained in the
 220 translation to empirical multiplexes, their three distributions (where three is equal to the number of
 221 multiplexed targets) from the singleplex reaction were compared with their corresponding distributions in
 222 empirical multiplex reactions. As illustrated in **Fig. 3a-c**, distributions of three different multiplex assays are
 223 visualised with their relative mean values represented by the dashed/dotted lines. The figures show the
 224 capabilities of the “c” parameter to maintain distance information going from simulation to empirical test.
 225 It can be observed that in most cases, the location of the parameter distribution for each target is
 226 maintained. In other situations, the distribution may be shifted from the singleplex events; however, the
 227 relative distance relationship of “c” values is kept. **Fig. 3a** illustrates the “c” parameter distribution of a low-
 228 rank ADS/MDS multiplex, showing overlaps for all the three singleplex assays in both simulated and empirical
 229 multiplexes. As distances among amplification curve shapes can significantly affect the ACA classifier,
 230 reduced performance is expected for multi-target identification. Another distribution trend among multiplex
 231 assays is represented in **Fig. 3b**, where the selected Primer Mix (PM3.01) has a high simulated ADS value
 232 (0.117) but low MDS (0.003). Moreover, we reported that the ADS value for distributions in **Fig. 3c** equals
 233 0.138, which differs only 0.21 from the combination PM3.01. However, PM3.12 has an MDS value of 0.075,
 234 representing an increase of 0.072 compared to PM3.01. This highlights the importance of considering
 235 minimum distances between “c” parameter distributions of the two closest targets: a small MDS value
 236 indicates a less separable group of target clusters, resulting in low ACA accuracies for multi-pathogen
 237 identification in a single fluorescent channel reaction. To numerically report how distributions are related in
 238 the translation from simulated to empirical multiplexes, we calculated the *Rooted Mean Squared Error*
 239 (*RMSE*) as follows:

$$240 \quad RMSE = \sqrt{\frac{(\mathbf{D}_s - \mathbf{D}_m)^T (\mathbf{D}_s - \mathbf{D}_m)}{N_d}}$$

241 where \mathbf{D}_s and \mathbf{D}_m are vectors for distances among targets in singleplex and multiplex, respectively. *RMSE*
 242 values of all the 3-plex combinations range from 0.003 to 0.050, which are negligible considering the range
 243 of the “c” parameters. The ADS, MDS and *RMSE* values for all the 3-plex combinations are reported in
 244 **Supplementary Table 4**. These results emphasise that distances between simulated and empirical multiplex
 245 share high similarity across different ranks, ensuring that our scoring system (based on ADS and MDS) is not
 246 affected whether in singleplex or multiplex environments.



247
 248 **Fig. 3 | Relative “c” parameter distributions of three different multiplex assays.** **a.** Primer Mix 3.07 (PM3.07) illustrates the “c”
 249 parameter distribution of a low-rank ADS/MDS multiplex; **b.** PM3.01 as an example of high ADS but low MDS multiplexes. **c.**
 250 Multiplex assay with high ADS and MDS with clearly separated distributions. For each subplot, the left graph shows the
 251 distributions of “c” parameters for the Simulated Multiplex. The right plot represents the corresponding distributions
 252 according to the empirical multiplex data. The vertical dashed lines correspond to the mean of the distribution computed for
 253 different targets. To quantitatively verify that the distances are maintained in transition from simulated to empirical
 254 multiplexes, the *RMSE* of distances is calculated and displayed on the graph title.

255 **Accuracy of all the possible combinations in 3-plex assays.** One of the aims of the Smart-Plexer is to improve
 256 the classification of multiplex assays, in our case, related to the ACA method. As demonstrated in the
 257 previous section, distances among amplification curves of empirical multiplex assays are similar to those
 258 generated in simulated multiplexes. Therefore, leveraging ADS and MDS, simulated multiplexes can be used
 259 to rank each combination and find the optimal assays with the largest inter-target distances for the ACA
 260 classifier. To further demonstrate that the ADS and MDS are crucial to improving multi-target identification
 261 in single well PCR reactions, we assess the classification performance of the ACA method by using 10-fold
 262 cross-validation and the k-Nearest Neighbors (k-NN) algorithm. **Fig. 4a** shows a 3-D graph where both ADS
 263 and MDS of the “c” parameters are correlated to the ACA accuracy. Accuracy percentages ranged from
 264 98.63% to 100% for each multiplex. The rainbow plane, which is fitted with linear regression on all the
 265 visualised data points, represents the gradient of the classification accuracy, showing an upward trend as
 266 ADS and MDS increase, which is consistent with our hypothesis that the ACA classification performs better
 267 with larger inter-target distances. Moreover, the plane on the left of **Fig. 4a** has a grey highlight zone called
 268 Vacuumed Area, where data points cannot fall inside as it is mathematically impossible to have an average
 269 distance value smaller than the minimum distance. We also defined another area called Forbidden Area, as
 270 visualised in the rotated 3-D plot on the right of **Fig. 4a**, where it is expected that no point will be founded,
 271 provided high values for ADS and MDS.

272 Both 3-D plots have circled points labelled as the top combination (TOP), bottom combination with lowest
 273 ADS (BOT ADS), bottom combination with lowest MDS (BOT MDS), and outlier combination (OUTLIER), with
 274 ACA classification accuracies of 99.9%, 99.89%, 98.06%, 99.01%, 99.82% and 99.87%, respectively. Although
 275 the overall classification performance for all the 27 combinations shows a high average of $99.51\% \pm 0.41\%$,

276 an increase of 1.84% is observed for the top ADS/MDS data point compared to the bottom one.
277 Furthermore, as depicted in **Fig. 4b-e**, by applying 2-D t-distributed stochastic neighbor embedding (t-SNE)⁴²
278 visualisation on curves generated by the top and bottom-ranked primer combinations, more condensed
279 target clusters and better separated inter-target boundaries can be seen for top-ranked assays. This results
280 in more distinguishable curve shapes and larger curve distances among targets, which benefits the ACA
281 classification. Numerical analysis of the visualised clusters was assessed using the Mean Silhouette Scores
282 (MSS). As reported by Kaufman et al. 2009, Silhouette scores between 0.51-0.70 are considered more
283 effective in cluster separation than values below 0.50⁴³. The reported MSS scores show significantly larger
284 inter-cluster distances for the top combinations, with values higher than 0.61 as opposed to the bottom
285 ones of less than 0.27 (in **Supplementary Table 4**, we also report ADS, MDS, MSS and ACA accuracies for
286 each combination of the 3-plex experiment). This finding proves that the ADS and MDS metrics are valid
287 indicators for predicting optimal primer set combinations for the ACA classifier. Relying on the Smart-Plexer
288 for selecting multiplex assays from singleplexes, the likelihood of accurate multi-target identification in a
289 single fluorescent channel reaction is significantly increased using the ACA methodology.

290 As mentioned above, **Fig. 4a** highlights the presence of outlier combinations where small ADS/MDS with
291 high ACA accuracy are reported (instead, low accuracy for the ACA classifier is expected). However, the
292 existence of such data points does not deny the effectiveness of the proposed method. It is important to
293 emphasise that the overall ACA accuracy for 3-plexes is inherently high because of the low levels of
294 multiplexing. Classifying three different curve shapes does not represent a major challenge for this Machine
295 Learning method, and targets with minor curve-shape differences can be easily separated in the feature
296 space. Considering this, along with the prevalent randomness that exists in the ACA method for 3-plex,
297 accuracies higher and lower than expected may occur in the given dataset. In fact, in the area with low
298 ADS/MDS, we can observe a large standard deviation for accuracies among data points which fall beneath
299 and above the fitted plane. Regardless of the accidentally high accuracies and low ADS/MDS caused by
300 randomness, **Figure 4f-g** evidence that these outlier combinations will face more challenges when used for
301 multi-target identification in larger scale multiplexes (or high-level multiplexing). In the outliers, the mapped
302 target clusters are largely overlapped with unclear boundaries and small MSS even in 3-plex assays.
303 Therefore, we will demonstrate in the next section that the higher the level of multiplexing is, the more
304 difficult the target separations are in the feature space when using these outliers.

305 Although low ADS/MDS combinations may occasionally show good performances, the proposed method
306 ensures that all predicted optimal multiplex assays with high ADS/MDS show high accuracies in ACA and
307 never the opposite. As illustrated in the 3-D plots of **Figure 4a**, the forbidden area (the red triangular prism)
308 has no data point falling in, which highlights the effectiveness of the ADS/MDS ranking system. This is a first
309 ever demonstration that multiplex assays tailored to the ACA method can be *in-silico* developed starting
310 from singleplex PCR reactions. This not only increases the likelihood of accurate multi-pathogen
311 identification, but also allows for a higher level of multiplexing in a single fluorescent channel. To
312 demonstrate the capabilities of the Smart-Plexer in developing optimal high-level multiplex assays for data-
313 driven approaches, in the following section, we assess its performance with seven different targets.

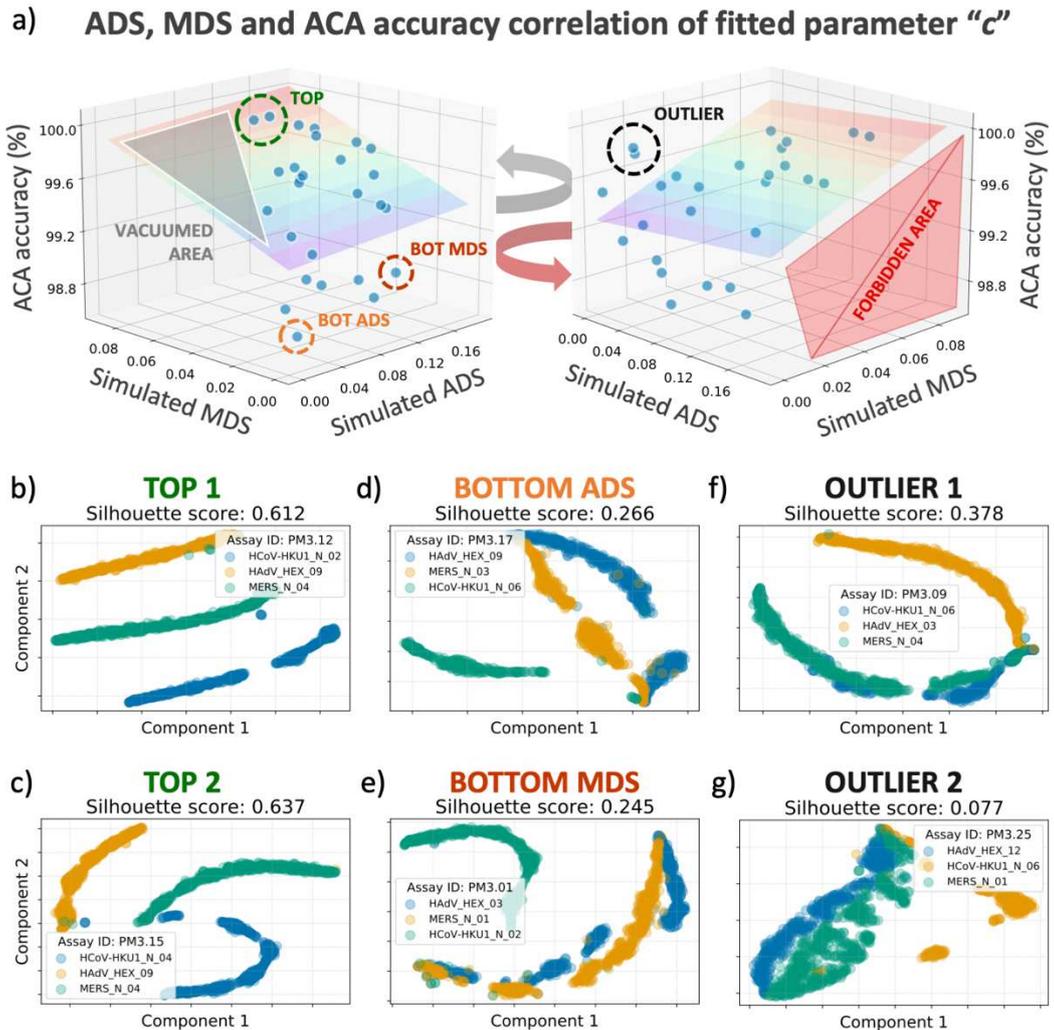


Fig. 4 | The influence of ADS/MDS on the ACA performance for all possible 3-plex combinations. a. 3-D plot of ACA classification accuracy for each combination versus simulated ADS and simulated MDS computed based on the “c” parameter. The rainbow plane is calculated using linear regression. In the left 3-D figure, the grey highlighted area is called Vacuumed Area, where simulated MDS is larger than simulated ADS (combinations in this area are mathematically impossible to be found). The right 3-D figure is a rotation of the left one, where a red is highlighted named Forbidden Area. In this region, high ADS/MDS combinations possess low ACA accuracies; however, no combinations were found. b-g. For the combination circled (TOP, BOT MDS, BOT ADS and OUTLIER) in a., 2-D t-SNE was applied on raw curves. In addition, for quantitative verification, the Mean Silhouette Scores (MSS) of target clusters were reported in the subplot title.

Smart-Plexer for development of 7-plex assays. In the previous section, our focus was on using a small number of targets to demonstrate that the developed ADS and MDS used to correlate distances between curves in both simulated and empirical multiplex assays were maintained. Moreover, accuracies among all the different combinations were evaluated using the ACA methodology, where high ADS/MDS multiplex assays show the highest likelihood of correct multi-target classification. These previous results indicate that the Smart-Plexer is a promising technique for optimal selection of primer set combinations in data-driven multiplexing.

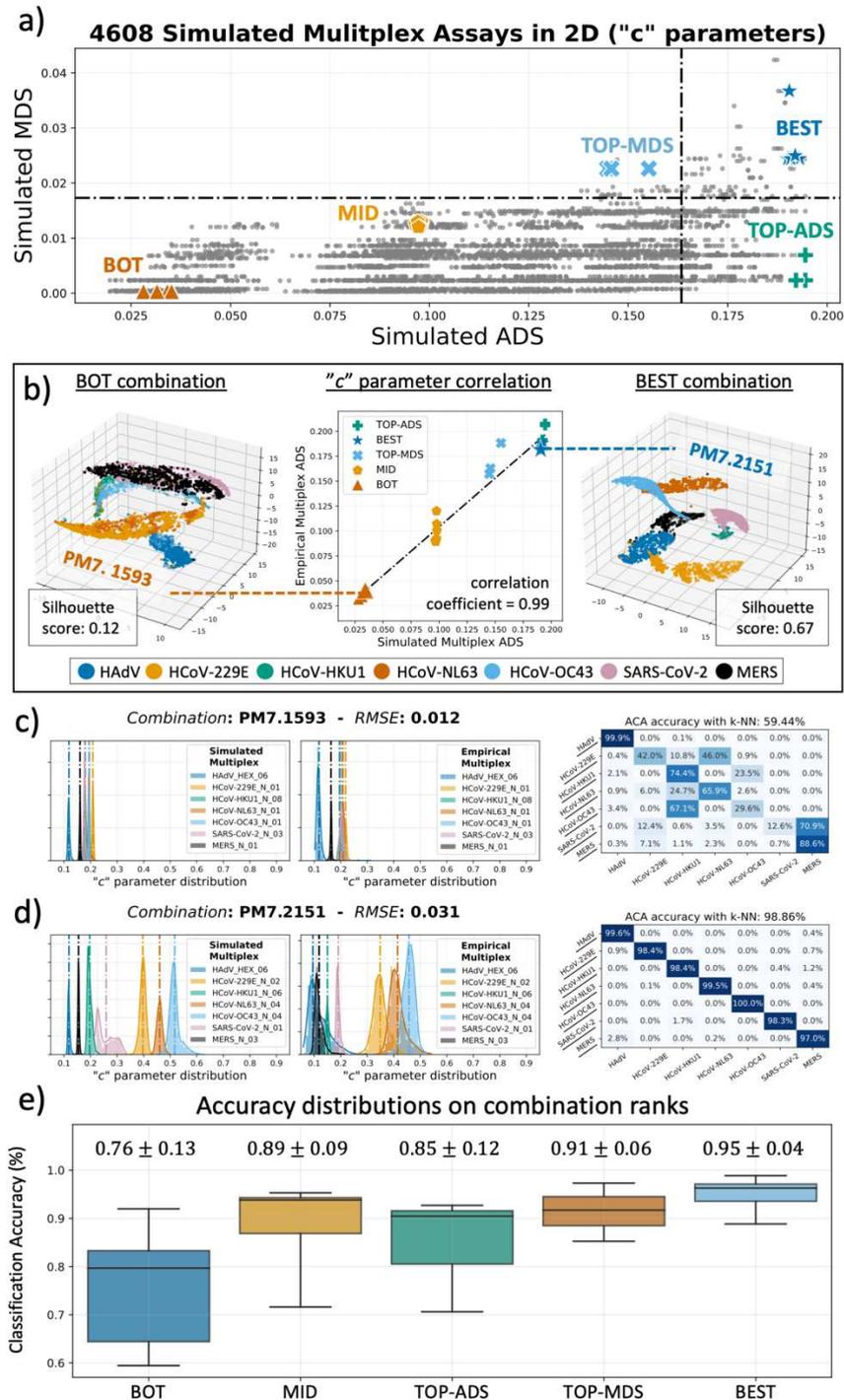
Next, we challenged the Smart-Plexer to develop an optimal 7-plex assay, which through the ACA method, is able to accurately identify the following Respiratory Tract Infection (RTI) pathogens in a single fluorescent channel using qdPCR: Human adenovirus (HAdV), Human coronavirus OC43 (HCoV-OC43), Human coronavirus HKU1 (HCoV-HKU1), Human coronavirus 229E (HCoV-229E), Human coronavirus NL63 (HCoV-NL63), Middle East respiratory syndrome-related coronavirus (MERS-CoV), and Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). We designed at least two different assays for each target, for a total of 24 singleplexes across the seven pathogens, as shown in **Supplementary Table 5**. Each primer set was

338 tested using synthetic DNA of its correspondent pathogenic target. Following the previous 3-plex
339 experimental workflow, the resulting raw curves were processed, fitted, and passed to the Smart-Plexer to
340 calculate all possible 7-plex combinations (N=4608) and compute their ADS/MDS. Based on “c” parameter
341 distances from fitted simulated multiplexes, **Figure 5a** shows how the ADS and MDS can be visualised in a
342 two-dimensional space. By considering the mean and standard deviation of the two scores, we set up
343 boundaries to the ADS/MDS distribution for all the combinations and divided the space into four separate
344 regions, with the purpose of showing how empirical multiplexes would perform for the ACA method
345 depending on their ADS/MDS. The black horizontal segmented line in **Figure 5a** divides high and low MDS,
346 and the vertical one separates the two ADS regions, resulting in four distinct areas. By testing different
347 multiplexes from each of these regions, we are aiming to further demonstrate that chance of developing a
348 reliable multiplex can vary based on the selected regions or selection criteria. Therefore, we chose multiplex
349 assays from different areas and categorised them into five classes, which were empirically tested with
350 synthetic DNA in qdPCR: BOT (N=6), MID (N=6), BEST (N=6), TOP-ADS and TOP-MDS (N=6) values (detailed
351 selection criteria are reported in the methodology section).

352 After the empirical testing, the distances of the “c” parameters of each selected multiplex were compared
353 to the simulated one, resulting in a correlation coefficient of 0.99, as shown in the middle graph of **Figure**
354 **5b**. Moreover, empirical multiplex amplification events were visualised using 3-D t-SNE, and distances across
355 target clusters were calculated with the MSS. As shown in the left plot of **Figure 5b**, clusters of the selected
356 BOT combination have an MSS of 0.12, whereas for the BEST one the score is 0.67. It can be observed that
357 there is a clear difference in clustering between the two selected multiplex assays, where the BEST one
358 shows clear separation among different targets (in line with the 3-plex results), and is expected to converge
359 in better ACA classification. The opposite scenario is shown in the BOT combination.

360 We validated that in higher level multiplexing, distance distributions of the “c” parameters were still
361 maintained from simulated to empirical testing; therefore, we computed the *RMSE* of the chosen tested
362 combinations. **Figure 5c-d** illustrate side-by-side “c” parameter distributions for each target in both
363 simulated (left) and empirical (right) multiplexes, showing a small *RMSE* for both BOT and BEST assays (0.012
364 and 0.031), and confirming the distance-maintaining hypothesis validated in the 3-plex experiments.
365 Moreover, we evaluate the ACA accuracy using training and testing datasets obtained in different
366 experimental settings (different days, operators, and reagents) to ensure the reproducibility of the
367 methodology. As expected, the performance of the BEST combination was significantly higher than the BOT
368 one, with a 39.42% increase in accuracy. Furthermore, in **Supplementary Table 6**, we reported the ADS, MDS
369 and accuracy values for the 24 selected multiplex assays. In **Supplementary Fig. 2**, we also visualised the
370 standard curve for each target using the BEST 7-plex assay to evaluate primer sensitivity and specificity. The
371 chosen multiplex reached a limit of quantification equal to 10^2 for all the respiratory pathogens using
372 synthetic DNA in real-time PCR.

373 As described before, ACA performances were evaluated using training and testing datasets from different
374 experimental settings with the same sample size. All the selected 24 multiplexes were empirically tested,
375 and their multi-target identification performances were assessed. In **Figure 5e**, accuracies and standard
376 deviations of each group of multiplexes were reported and visualised as box plots. The best-combination
377 group scored an average (\pm standard deviation) classification performance of 95% (\pm 0.04%) using a k-NN
378 classifier, which is the highest average and the lowest standard deviation among all the groups. There is a
379 decreasing trend in the average accuracy, and an increasing trend in the standard deviation as the ADS/MDS
380 values become smaller. Previously, the 3-plex validation showed the presence of outliers in low ADS/MDS
381 rank with high ACA classification accuracy, which is also observed in these 7-plex tests. However, the
382 standard deviation indicates that the Smart-Plexer does provide a robust and solid solution (even at high-
383 level multiplexing) to significantly increase the likelihood of choosing an optimal multiplex for data-driven
384 multiplexing (i.e. ACA methodology).



385
386
387
388
389
390
391
392
393
394
395
396
397
398

Fig. 5 | Validation of Smart-Plexer based on 7-plex assays. **a.** 2-D ranking results for all 4608 combinations in 7-plex based on simulated ADS and simulated MDS. The plot is divided into four regions to explore the relationship between ACA performance and ADS/MDS. Combinations from five different classes (BOT, MID, BEST, TOP-ADS, TOP-MDS) were selected for multiplex empirical testing. **b.** The 2-D plot in the middle depicts the relationship between empirical and simulated scores based on "c" parameters. Enlarged data points for one of the BOT (PM7.1593) and BEST (PM7.2151) combinations are visualised with 3-D t-SNE on raw curves, and the corresponding Silhouette scores are calculated. **c-d.** Simulated and Empirical "c" distributions of the selected combinations (PM7.1593 and PM7.2151) are plotted (RMSE values in subplot titles). The vertical dashed lines correspond to the mean of the distribution computed for different targets. On the right, the confusion matrixes of ACA performance for both cases are presented, and overall accuracy using k-NN is reported in the title. True labels are on the y-axis and ACA predicted labels are on the x-axis (each target sensitivity is also reported in percentage). **e.** The box plot of ACA classification accuracy for each selected group. The mean and standard deviation of ACA accuracy on empirical multiplexes are calculated and shown on each box bar.

399 **Clinical validation results.** The final step was to validate that the Smart-Plexer is capable of easing the
400 laboratory workload in developing multiplex assays. After testing six potential best combinations based on
401 ADS/MDS, we selected the one with the highest ACA classification accuracy on synthetic DNA (PM7.2151).
402 To clinically validate the selected 7-plex for multi-pathogen identification, inactivated clinical samples were
403 purchased from Randox Laboratories (UK) and extracted using a gold standard kit (QIAGEN mini amp). The
404 extracted samples were used as the testing dataset (7,638 positive amplification reactions), while curves
405 resulting from synthetic DNA amplification reactions (5,207 positive amplification reactions) were the
406 training. The classifier used was a k-NN with the number of neighbours equal to 10. As shown in **Table 1**, a
407 total of 14 positive samples were classified in qdPCR using the ACA methodology. The predicated label of a
408 sample is given by selecting the most predictable label within all the in-sample curves. The confidence level
409 was given as the percentage of the amplification curves with the most predicted label. We correctly
410 identified all pathogens using the Smart-Plexer selected candidate assay, where most of the targets were
411 recognised with high confidence (median = 95.46%).

412 It is important to note that in this study we faced a seven-class classification problem, where the accuracy
413 of a “random guess” (or a random classifier as convention) equals 14.3% under a balanced dataset. All the
414 confidence levels were much higher than the random guess accuracy, indicating solid and robust predictions
415 with the selected optimal multiplex assay. Although the number of clinical samples was limited by the
416 number of pathogens provided by the manufacturer, the proposed framework, in combination with the ACA
417 methodology, achieved a highly accurate identification of multiple pathogens by using an optimal multiplex
418 assay in a single fluorescent channel reaction. The Smart-Plexer can leverage the capability of the data-
419 driven multiplexing to an easy-to-develop, robust, and cost-effective molecular diagnostic solution.

420 Discussion

421 In this work, we developed the Smart-Plexer, an innovative framework which combines wet-lab
422 experiments and computational algorithms to generate optimal multiplex assays for data-driven approaches
423 using real-time PCR data. The method leverages mathematical metrics to construct an advanced ranking
424 system to increase the throughput of conventional molecular tests by optimising their chemical peculiarities.
425 To reveal the potential of this powerful approach, we demonstrated it with a recently reported machine
426 learning method, named Amplification Curve Analysis (ACA), which is capable of identifying multiple nucleic
427 acid targets in a single fluorescent channel with conventional PCR instruments. As the ACA leverages kinetic
428 information encoded in the amplification curve, multiple targets can be classified based on the unique
429 interaction with their assigned primer sets. However, constructing different amplification curve shapes for
430 each multiplexed target is one of the major challenges for the ACA approach. The Smart-Plexer solves this
431 problem by providing an easy-to-use framework for multiplex assay development, enabling high-level and
432 highly accurate data-driven multiplexing.

433 This study shows the progression of the Smart-Plexer starting from a simple three-target classification
434 problem. From the wet-lab testing of three singleplex assays for each of the three targets, a total of 27
435 combinations (in our case 3-plex assays) can be generated *in silico* (simulated multiplex) and ranked based
436 on the mathematical curve-shape distances. Using synthetic DNA in qdPCR and a single fluorescent channel,
437 the assays were empirically tested (empirical multiplex), and the ACA classification accuracies were
438 evaluated for all the possible combinations. The distance scores computed from the Smart-Plexer for
439 multiplex assay ranking are linearly correlated between simulated and empirical multiplexes. Moreover, we
440 showed a further correlation between high-rank multiplexes and a high probability of increasing the ACA
441 accuracies, confirming that the metrics used in this novel framework are theoretically connected to the
442 distance measurement of the machine learning classifier.

443 As the complexity of developing multiplex assays exponentially increases with the number of targets, we
444 further challenged the Smart-Plexer by designing a 7-plex assay to identify common respiratory tract
445 infection (RTI) pathogens. Consistent with the 3-plex validation, the correlation between simulated and
446 empirical multiplex is also maintained in 7-plex. Regarding the ACA classification, it is logical that higher
447 similarities among curves exist in a scenario with a higher number of targets, making it harder to develop
448 multiplex assays. Nevertheless, the Smart-Plexer brilliantly generated an optimal multiplex assay, which
449 correctly identified pathogens presented in 14 commercial clinical samples. It was further demonstrated
450 that, since ACA is a clustering method, it requires a large minimum distance between the two closest clusters

451 and a large average distance among all clusters in the multiplex. Therefore, the Smart-Plexer ranking system
452 enables the development of optimal multiplex assays for data-driven multiplexing.

453 Apart from the scalability of multiplexing that the Smart-Plexer can provide to the ACA method, we
454 demonstrated for the first time that machine learning approaches can be applied to probe-based
455 multiplexes, in our case, TaqMan. Probe-based assays, together with the use of intercalating dyes and
456 isothermal chemistries, are expanding the boundaries of data-driven multiplexing and opening new windows
457 for its application in commercial, research and clinical fields. The Smart-Plexer eases the development of any
458 novel multiplex panel or molecular assays, enabling the use of the ACA as an emerging diagnostic tool.
459 Through this hybrid method, it is possible to select the highest rank combination *in silico* with wet-lab tested
460 singleplexes, avoiding performing expensive and time-consuming multiplex assay development phases.

461 While this novel framework is validated with high-level multiplexing (7-plex), it is essential to highlight that
462 distances between amplification curves can be a limiting factor in single fluorescent channel multiplexing.
463 This affects the Smart-Plexer since the inter-target differences of fitting parameters considered for the
464 distance measurement become smaller as the target number increases. In this work, we use linear distance
465 measurements, but more advanced metrics (e.g. Minkowski, Chebyshev or Cosine) can be adopted to
466 improve the ranking performance. Moreover, when a higher level of multiplexing is required, the use of
467 probe-based chemistries such as TaqMan comes handy. By leveraging the optical capability of real-time PCR
468 instruments, a multiplex assay using multiple-channel detection can double or triple the number of targets
469 in a single reaction. All these strategies aim to improve the ACA classification through a more innovative
470 development from the chemistry perspective, while from the machine learning view, the current classifiers
471 rely on state-of-the-art algorithms which shine for their robustness but are limited for tailoring to specific
472 datasets. We previously demonstrated that more advanced classifiers such as convolutional neural networks
473 (CNN) could extend the ACA capability to classify targets for higher-level multiplex assays. However, as a
474 novel technique, data-driven multiplexing requires more optimisation and development of algorithms.

475 The Smart-Plexer represents a solution for developing multiplex assays by utilising both empirical testing
476 and *in-silico* computation. The hybrid nature of this framework still requires wet-lab experiments; therefore,
477 certain limitations exist in terms of staff training and time requirements. However, our future work will focus
478 on the full automation of developing such assays. Novel methodologies to predict amplification curve
479 behaviours will be developed. One example is the brand-new algorithm for designing multiplex PCR primers
480 using Dimer Likelihood Estimation by Xie et al. 2021²¹. Another future aspect of this research is to further
481 increase inter-target curve shape differences. An example would be in probe-based chemistries, where
482 modifying amplification curve shapes can be achieved by changing the concentration levels of the
483 fluorescent probe. In this way, we can enlarge inter-target distances of amplification curves to ease the ACA
484 classification with better clustering performance. All the above-mentioned future works will inspire the use
485 of the ACA method for a broad range of applications and significantly increase its flexibility and scalability.

486 In this work, we present for the first time a complete pipeline for developing optimal multiplex assays and
487 open the usage of the ACA method to the broad scientific community. We finally unlock the development
488 of novel molecular tests that will not only enable simpler molecular diagnostics but also make them
489 affordable and available to everyone.

490 References

- 491 1. Heid, C. A., Stevens, J., Livak, K. J. & Williams, P. M. Real time quantitative PCR. *Genome Res.* **6**, 986–
492 994 (1996).
- 493 2. Rebrikov, D. V. & Trofimov, D. Yu. Real-time PCR: A review of approaches to data analysis. *Appl*
494 *Biochem Microbiol* **42**, 455–463 (2006).
- 495 3. Higuchi, R., Fockler, C., Dollinger, G. & Watson, R. Kinetic PCR Analysis: Real-time Monitoring of DNA
496 Amplification Reactions. *Nat Biotechnol* **11**, 1026–1030 (1993).
- 497 4. Velavan, T. P. & Meyer, C. G. COVID-19: A PCR-defined pandemic. *International Journal of Infectious*
498 *Diseases* **103**, 278–279 (2021).
- 499 5. Wölfel, R. *et al.* Virological assessment of hospitalized patients with COVID-2019. *Nature* **581**, 465–
500 469 (2020).

- 501 6. Tan, C. *et al.* Applications of digital PCR in COVID-19 pandemic. *VIEW* **2**, 20200082 (2021).
- 502 7. Collier, D. A. *et al.* Point of Care Nucleic Acid Testing for SARS-CoV-2 in Hospitalized Patients: A
503 Clinical Validation Trial and Implementation Study. *Cell Reports Medicine* **1**, 100062 (2020).
- 504 8. Rodriguez-Manzano, J. *et al.* Handheld Point-of-Care System for Rapid Detection of SARS-CoV-2
505 Extracted RNA in under 20 min. *ACS Cent. Sci.* **7**, 307–317 (2021).
- 506 9. Vandenberg, O., Martiny, D., Rochas, O., van Belkum, A. & Kozlakidis, Z. Considerations for
507 diagnostic COVID-19 tests. *Nat Rev Microbiol* **19**, 171–183 (2021).
- 508 10. Scohy, A. *et al.* Low performance of rapid antigen detection test as frontline testing for COVID-19
509 diagnosis. *Journal of Clinical Virology* **129**, 104455 (2020).
- 510 11. Sreenath, K. *et al.* Coinfections with Other Respiratory Pathogens among Patients with COVID-19.
511 *Microbiology Spectrum* **9**, e00163-21 (2021).
- 512 12. Khaddour, K., Sikora, A., Tahir, N., Nepomuceno, D. & Huang, T. Case Report: The Importance of
513 Novel Coronavirus Disease (COVID-19) and Coinfection with Other Respiratory Pathogens in the Current
514 Pandemic. *Am J Trop Med Hyg* **102**, 1208–1209 (2020).
- 515 13. Zhu, H. *et al.* The vision of point-of-care PCR tests for the COVID-19 pandemic and beyond. *TrAC*
516 *Trends in Analytical Chemistry* **130**, 115984 (2020).
- 517 14. Mahony, J. B. *et al.* Cost Analysis of Multiplex PCR Testing for Diagnosing Respiratory Virus
518 Infections. *Journal of Clinical Microbiology* **47**, 2812–2817 (2009).
- 519 15. Barenfanger, J., Drake, C., Leon, N., Mueller, T. & Troutt, T. Clinical and Financial Benefits of Rapid
520 Detection of Respiratory Viruses: an Outcomes Study. *Journal of Clinical Microbiology* **38**, 2824–2828 (2000).
- 521 16. Khodakov, D., Wang, C. & Zhang, D. Y. Diagnostics based on nucleic acid sequence variant profiling:
522 PCR, hybridization, and NGS approaches. *Advanced Drug Delivery Reviews* **105**, 3–19 (2016).
- 523 17. Edwards, M. C. & Gibbs, R. A. Multiplex PCR: advantages, development, and applications. *Genome*
524 *Res.* **3**, S65–S75 (1994).
- 525 18. Caliendo, A. M. Multiplex PCR and Emerging Technologies for the Detection of Respiratory
526 Pathogens. *Clinical Infectious Diseases* **52**, S326–S330 (2011).
- 527 19. van der Zee, A. *et al.* Multi-centre evaluation of real-time multiplex PCR for detection of
528 carbapenemase genes OXA-48, VIM, IMP, NDM and KPC. *BMC Infect Dis* **14**, 27 (2014).
- 529 20. Rodriguez-Manzano, J. *et al.* Simultaneous Single-Channel Multiplexing and Quantification of
530 Carbapenem-Resistant Genes Using Multidimensional Standard Curves. *Anal. Chem.* **91**, 2013–2020 (2019).
- 531 21. Xie, N. G. *et al.* Designing highly multiplex PCR primer sets with Simulated Annealing Design using
532 Dimer Likelihood Estimation (SADDLE). *Nat Commun* **13**, 1881 (2022).
- 533 22. Meuzelaar, L. S., Lancaster, O., Pasche, J. P., Kopal, G. & Brookes, A. J. MegaPlex PCR: a strategy for
534 multiplex amplification. *Nat Methods* **4**, 835–837 (2007).
- 535 23. Farrar, J. S. & Wittwer, C. T. Chapter 6 - High-Resolution Melting Curve Analysis for Molecular
536 Diagnostics. in *Molecular Diagnostics (Third Edition)* (ed. Patrinos, G. P.) 79–102 (Academic Press, 2017).
537 doi:10.1016/B978-0-12-802971-8.00006-7.
- 538 24. Zhang, Q., Yang, F., Gao, J., Zhang, W. & Xu, X. Development of multiplex TaqMan qPCR for
539 simultaneous detection and differentiation of eight common swine viral and bacterial pathogens. *Braz J*
540 *Microbiol* **53**, 359–368 (2022).
- 541 25. Lee, J. *et al.* Identification of *Lactobacillus sakei* and *Lactobacillus curvatus* by multiplex PCR-based
542 restriction enzyme analysis. *Journal of Microbiological Methods* **59**, 1–6 (2004).
- 543 26. Moniri, A. *et al.* Framework for DNA Quantification and Outlier Detection Using Multidimensional
544 Standard Curves. *Anal. Chem.* **91**, 7426–7434 (2019).
- 545 27. Rawson, T. M., Peiffer-Smadja, N. & Holmes, A. Artificial Intelligence in Infectious Diseases. in
546 *Artificial Intelligence in Medicine* (eds. Lidströmer, N. & Ashrafian, H.) 1–14 (Springer International
547 Publishing, 2020). doi:10.1007/978-3-030-58080-3_103-1.
- 548 28. Ozkok, F. O. & Celik, M. A hybrid CNN-LSTM model for high resolution melting curve classification.
549 *Biomedical Signal Processing and Control* **71**, 103168 (2022).

550 29. Jacky, L. *et al.* Robust Multichannel Encoding for Highly Multiplexed Quantitative PCR. *Anal. Chem.*
551 **93**, 4208–4216 (2021).

552 30. Moniri, A. *et al.* Amplification Curve Analysis: Data-Driven Multiplexing Using Real-Time Digital PCR.
553 *Anal. Chem.* **92**, 13134–13143 (2020).

554 31. Moniri, A., Miglietta, L., Holmes, A., Georgiou, P. & Rodriguez-Manzano, J. High-Level Multiplexing
555 in Digital PCR with Intercalating Dyes by Coupling Real-Time Kinetics and Melting Curve Analysis. *Anal. Chem.*
556 **92**, 14181–14188 (2020).

557 32. Miglietta, L. *et al.* Coupling Machine Learning and High Throughput Multiplex Digital PCR Enables
558 Accurate Detection of Carbapenem-Resistant Genes in Clinical Isolates. *Frontiers in Molecular Biosciences* **8**,
559 (2021).

560 33. Malpartida-Cardenas, K. *et al.* Single-channel digital LAMP multiplexing using amplification curve
561 analysis. *Sensors & Diagnostics* **1**, 465–468 (2022).

562 34. Elnifro, E. M., Ashshi, A. M., Cooper, R. J. & Klapper, P. E. Multiplex PCR: Optimization and
563 Application in Diagnostic Virology. *Clin Microbiol Rev* **13**, 559–570 (2000).

564 35. Markoulatos, P., Siafakas, N. & Moncany, M. Multiplex polymerase chain reaction: A practical
565 approach. *Journal of Clinical Laboratory Analysis* **16**, 47–51 (2002).

566 36. Rachlin, J., Ding, C., Cantor, C. & Kasif, S. Computational tradeoffs in multiplex PCR assay design for
567 SNP genotyping. *BMC Genomics* **6**, 102 (2005).

568 37. Ozaki, Y. *et al.* Cost-efficient multiplex PCR for routine genotyping of up to nine classical HLA loci in
569 a single analytical run of multiple samples by next generation sequencing. *BMC Genomics* **16**, 318 (2015).

570 38. Miglietta, L. *et al.* An adaptive filtering framework for non-specific and inefficient reactions in
571 multiplex digital PCR based on sigmoidal trends. *bioRxiv* 2022.04.11.487847 (2022)
572 doi:10.1101/2022.04.11.487847.

573 39. Spiess, A.-N., Feig, C. & Ritz, C. Highly accurate sigmoidal fitting of real-time PCR data by introducing
574 a parameter for asymmetry. *BMC Bioinformatics* **9**, 221 (2008).

575 40. Liu, W. *et al.* A Novel Sigmoid-Function-Based Adaptive Weighted Particle Swarm Optimizer. *IEEE*
576 *Transactions on Cybernetics* **51**, 1085–1093 (2021).

577 41. Ukalska, J. & Jastrzębowski, S. Sigmoid growth curves, a new approach to study the dynamics of the
578 epicotyl emergence of oak. *Folia Forestalia Polonica* **61**, 30–41 (2019).

579 42. van der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *Journal of Machine Learning Research*
580 **9**, 2579–2605 (2008).

581 43. Partitioning Around Medoids (Program PAM). in *Finding Groups in Data* 68–125 (John Wiley & Sons,
582 Ltd, 1990). doi:10.1002/9780470316801.ch2.

583

584 Methods

585 **Synthetic double-stranded DNA templates.** Double-stranded synthetic DNA was used in this study to develop
586 and assess the performance of all singleplex assays. In particular, we used the entire coding sequence of the
587 hexon protein gene (HEX gene) for human adenovirus (HAdV), and the nucleocapsid protein gene (N gene)
588 of human coronavirus OC43 (HCoV-OC43), HKU1 (HCoV-HKU1), 229E (HCoV-229E), NL63 (HCoV-NL63),
589 Middle East respiratory syndrome-related coronavirus (MERS-CoV) and severe acute respiratory syndrome
590 coronavirus 2 (SARS-CoV-2). The following NCBI accession numbers were used as references for the gBlock
591 synthesis: NC_001405, NC_006213, NC_006577, NC_002645, NC_005831, NC_019843 and NC_045512,
592 respectively. The synthetic constructs were used for qPCR experiments when determining the limit-of-
593 quantification of each PCR assay, and in qdPCR experiments for generating the dataset used in the simulation
594 of the multiplexes and their empirical testing. The gene fragments (ranging from 1,134 to 1,558 bp) were
595 purchased from Integrated DNA Technologies Ltd. (IDT) and resuspended in Tris-EDTA buffer to 10 ng/ μ l
596 stock solutions (stored at -80°C until further use). The concentrations of all DNA stock solutions were
597 determined using a Qubit 3.0 fluorimeter (Life Technologies).

598

599 **Clinical samples.** Whole pathogen control panels were purchased from Randox Laboratories Ltd, including
600 MERS-CoV (catalog no. QAV154181), CoV-OC43, NL63 (catalog no. QAV164189), and SARS-CoV-2 (catalog
601 no. SCV2QC). Samples were extracted using the QIAamp Viral RNA Mini Kits (catalog no. 52906). Viral nucleic
602 acid was extracted using the manufacturer-recommended protocol⁴⁴. Viral RNA was reverse transcribed to
603 cDNA using Fluidigm reverse transcription master mix (catalog no. SKU 100-6299). Viral cDNA was further
604 pre-amplified using Fluidigm Preamp master mix (catalog no. PN 100-5744). Reverse transcription and pre-
605 amplification were conducted according to the Fluidigm manufacturer's protocol (Fluidigm document
606 number: 101-7571 A2 and 100-5876 C2).

607

608 **PCR assay design.** The sequences of each gene were downloaded from the GenBank website⁴⁵. Based on the
609 comprehensive analyses and alignments of each type using the MUSCLE algorithm⁴⁶, primers were
610 specifically designed to amplify all sequence variations within each gene belonging to their specific target
611 (inclusivity) and to exclude closely related but not inclusive sequences (exclusivity). Design and in-silico
612 analysis were conducted using GENEious Prime 2022.0.1⁴⁷. Primer characteristics were analysed through
613 IDT OligoAnalyzer software⁴⁸ using the J. SantaLucia thermodynamic table for melting temperature (T_m)
614 evaluation, hairpin, self-dimer, and cross-primer formation⁴⁹. To confirm the specificity of the real-time
615 digital PCR assays, the primers were first evaluated in a singleplex PCR environment to address their
616 specificity and sensitivity for both singleplex and multiplex assays. All primers were synthesised by IDT
617 (Coralville, IA, United States). Details on singleplex and multiplex assays are provided in **Supplementary Table**
618 **7 and 10**, along with the primer sequences for both 3-plex and 7-plex.

619

620 **Real-time digital PCR.** For real-time amplification experiments, we used the BioMark HD (Fluidigm) and the
621 QIAquant 96 5plex (catalog no. 9003011). The master mix used was the PrimeTime Gene Expression Master
622 Mix from Integrated DNA Technologies (IDT, catalog no. 1055772) supplemented with ROX passive reference
623 dye and pre-mixed following manufacture guidelines. The qdPCR was performed with Fluidigm qdPCR 37k
624 integrated fluidic circuits (IFC) (catalog no. SKU100-6152) and was supplemented with Fluidigm 20X GE
625 loading buffer (PN 85000746). The priming and loading steps of the IFC were followed as the supplier's
626 protocol (Fluidigm document number: 100-6896 Rev 03). Each amplification mix for the qdPCR experiment
627 contained 3 μ l 2X IDT PrimeTime Gene Expression Master Mix (with passive ROX), 0.6 μ l 20X GE, 0.6 μ l 10X
628 Primer mixture, 1.8 μ l DNA templates from synthetic DNA, pre-amplified cDNA, or controls, and to bring the
629 final volume to 6 μ l. A total of 4.5 μ l of reaction mix was transferred to each inlet (or panel) of a Fluidigm
630 37k IFC for the thermal cycling step. Thermal-cycle conditions consisted of a hot start step for 3 min at 95
631 $^{\circ}\text{C}$, followed by 45 cycles at 95 $^{\circ}\text{C}$ for 15 seconds and 60 $^{\circ}\text{C}$ for 45 seconds. Real-time data of the amplification
632 events were exported as a text file for each bulk by Fluidigm Digital PCR Analysis software (version 4.1.2)

633

634 **Limit-of-quantification.** We used real-time PCR from QIAGEN (QIAQuanta96) to evaluate the Limit-of-
 635 quantification (LoQ) of the selected 7-plex assay. Standard curves were generated with synthetic DNA
 636 ranging from 10^7 to 10^1 , apart from SARS-CoV-2 whose concentration was from 10^5 to 10^1 because of
 637 limitations due to pandemic suppliers (IDT). PCR data were extracted and processed according to the data
 638 processing step. Standard curve plots and statistical values are reported in **Supplementary Fig. 2**. The
 639 Absence of amplification signals was detected in Negative Template Control (NTC)

640

641 **Data processing.** The processing of raw amplification curves is comprised of three parts. Firstly, to ensure all
 642 curves start from approximately zero fluorescence value and to normalise the starting cycles of the curve
 643 across the entire time series, the background information was removed, which can be expressed as:

644

$$Fl_{br}(t) = Fl(t) - avg_{back}$$

645 where $Fl_{br}(t)$ represents a curve with the background removed and $Fl(t)$ is the raw fluorescence values
 646 for each cycle $t = 1, 2, \dots, T$. Here T indicates the total number of cycles for each amplification curve (45 in
 647 our case), and avg_{back} is the average background value. In order to avoid instrumental noise commonly
 648 found at the beginning of the PCR reaction, the avg_{back} value was estimated as the average value of the
 649 first several cycles' fluorescence, excluding the initial ones. In our case, five cycles were considered for the
 650 flat phase and the first three cycles were skipped. Secondly, late amplification filtering was applied to select
 651 curves that reached the plateau phase. The basic idea is to estimate the cycle threshold value (Est_{Ct}) for
 652 each curve, which can be represented as:

653

$$Est_{Ct} = \min t$$

654

$$s.t. \frac{Fl_{br}(t) - F_{min}}{F_{max} - F_{min}} \geq F_{th}$$

655 where $t \in \{1, 2, \dots, T\}$, and F_{max} and F_{min} represent maximum and minimum fluorescence values of the
 656 entire reaction respectively for each curve. F_{th} is the fluorescence threshold and curves whose Est_{Ct} are
 657 above the cycle threshold ($C_t = 30$ as suggested by the manufacturer) were removed. Lastly, a filter was
 658 applied to remove non-sigmoidal curves with excessive noisy signals. The sigmoidal trend of a noisy curve
 659 may contain certain notches. Based on this feature, we estimated the first derivative of each curve using:

660

$$Fl_{br}'(t) = Fl_{br}(t) - Fl_{br}(t - 1), \quad t = 2, \dots, T$$

661 The number of zero-crossing points in $Fl_{br}'(t)$ is related to the number of notches in the curve. Therefore,
 662 noisy curves should have significantly more zero-crossing points in their first derivatives compared with
 663 smooth sigmoidal curves. The curves that satisfied the following condition were regarded as noisy and
 664 removed:

665

$$\sum_t \frac{-sgn[Fl_{br}'(t)] + 1}{2} > N_{zc}$$

666 where $sgn[\cdot]$ is the sign function and N_{zc} is the given threshold value ($N_{zc} = 9$ in our research).

667

668 **Five-parametric sigmoidal fitting.** Since amplification curves contain information such as background,
 669 plateau phase, and slope, we derived the most representative features of it using the sigmoidal equation.
 670 The chosen model in this study for curve fitting is the five-parametric sigmoid function, whose equation is
 671 given below:

672

$$f(t, \mathbf{p}) = \frac{a}{(1 + \exp^{-c(t-d)})^e} + b$$

673

$$\mathbf{p} = [a, b, c, d, e]^T$$

674 where t is the amplification cycle, \mathbf{p} is the parameter vector, $f(t, \mathbf{p})$ is the fluorescence at cycle t . The
 675 mathematical function of these parameters and their corresponding representations in amplification curves
 676 are shown in **Table 2** below:

677

678 To reduce optimisation iterations and unsuccessful fitting, we applied a pivot fitting on a subset of data (\mathbf{D}_s)
 679 to evaluate the optimal initial parameters \mathbf{p}_0^{opt} for the equation before searching on the entire dataset (\mathbf{D}).
 680 First, we defined a non-linear Least Square function $LS(\mathbf{p})$, whose equation is shown below:

$$681 \quad LS(\mathbf{p}) = \sum_{t=1}^T (f(t, \mathbf{p}) - Fl_{br}(t))^2$$

682 To apply the pivot fitting, we first initialised $\mathbf{p}_0 = [0,0,0,0,0]^T$. Then, for the i^{th} curve Fl_{br}^i within the dataset
 683 \mathbf{D}_s , the following optimisation problem was solved to find the fitted parameter vector:

$$684 \quad \mathbf{p}_i = \underset{\mathbf{B}_{low} < \mathbf{p} < \mathbf{B}_{up}}{\operatorname{argmin}} LS(\mathbf{p})$$

685 where the lower bound \mathbf{B}_{low} and the upper bound \mathbf{B}_{up} for all the parameters are -100 and 100,
 686 respectively. After all the curves were fitted, the mean vector of all the \mathbf{p}_i was used as the optimal \mathbf{p}_0^{opt} .

687 With the outcome from the pivot fitting, we fitted all curves in \mathbf{D} starting from \mathbf{p}_0^{opt} . In addition, to get
 688 better fitting performance, we increased the maximum number of fitting iterations (maxfev) to a sufficiently
 689 large value (1,000,000 in our case). The same \mathbf{B}_{low} and \mathbf{B}_{up} were used for the pivot fitting.

690

691 **Calculating Average Distance Score (ADS) and Minimum Distance Score (MDS) for multiplex assays.** There
 692 are four curve representations for calculating ADS and MDS, which are: raw curves (45-D), normalised curves
 693 (45-D), fitted parameters (5-D) and c parameter (1-D). Two steps were taken before the score calculation:
 694 (i) Extract the median feature vectors of each target for 45-D, 5-D and 1-D feature arrays. The median value
 695 was taken on each dimension, and the median feature vector with the same dimension was generated. It is
 696 assumed that the distribution of each target is Gaussian. However, outliers can affect the distribution
 697 unexpectedly. Therefore, the median value is a more robust representative compared to the average value,
 698 and \mathbf{N}_t median vectors corresponding to \mathbf{N}_t targets were constructed. (ii) Calculate Euclidean distance
 699 between each pair of targets, where given \mathbf{N}_t targets, the total number of distances \mathbf{N}_d is:

$$700 \quad \mathbf{N}_d = \binom{\mathbf{N}_t}{2} = \frac{\mathbf{N}_t(\mathbf{N}_t - 1)}{2}$$

701 The vector of distances for each pair of targets is defined as:

$$702 \quad \mathbf{S}_D = [d_{ij} \mid \text{for each } i = 2, \dots, \mathbf{N}_t, \quad j = 1, 2, \dots, i - 1]$$

703 where d_{ij} represents the Euclidean distance between extracted median vectors of target i and target j .
 704 With the constructed distance set, the ADS and MDS were calculated as the average and the minimum value
 705 of all elements in \mathbf{S}_D , respectively:

$$706 \quad ADS = \operatorname{mean}(\mathbf{S}_D)$$

$$707 \quad MDS = \operatorname{min}(\mathbf{S}_D)$$

708

709 **ACA methodology.** The Amplification Curve Analysis (ACA) methodology was developed by Moniri et al. in
 710 2020³⁰. For the first time, shapes of amplification curves from real-time PCR data were used for multiple
 711 target identification in a single fluorescent channel reaction, utilising data-driven algorithms. The ACA takes
 712 the entire amplification time series as input and uses machine learning to classify curves into different
 713 categories of targets. This approach highlights the significance of the kinetic information embedded in
 714 amplification curves. As previously reported, several classical machine learning methods (e.g. k-NN, Random
 715 Forest, Support Vector Machine) as well as deep-learning based approaches (e.g. Convolutional Neural
 716 Networks) can be applied to the time series³¹. In this article, a k-NN classifier with 10 neighbours was used
 717 for the ACA performance evaluation.

718

719 **Ranking system.** The inputs of the ranking system are simulated ADS and MDS. To increase the likelihood of
 720 choosing an optimal assay for data-driven multiplexing approaches, we considered assays with the highest
 721 ADS and MDS (S_{BEST}) selected from the entire combination set (S_{ALL}). Provided the number of the best

722 combinations to be selected as N_{BEST} and the number of total combinations as N_c , the following steps were
 723 applied:
 724

Algorithm 1

```

1:   Initialise  $N_{BEST}$  as required,  $S_{BEST} \leftarrow \emptyset$ 
2:   for  $n_e = N_{BEST}, N_{BEST} + 1, \dots, N_c$  do
3:      $S_{BEST}^{MDS} \triangleq \{x \mid x \text{ are the top } n_e \text{ combinations in } S_{ALL} \text{ with largest MDS}\}$ 
4:      $S_{BEST}^{ADS} \triangleq \{x \mid x \text{ are the top } n_e \text{ combinations in } S_{ALL} \text{ with largest ADS}\}$ 
5:      $S_{BEST} \leftarrow (S_{BEST}^{MDS} \cap S_{BEST}^{ADS}) \cup S_{BEST}$ 
6:     if  $|S_{BEST}| \geq N_{BEST}$ 
7:       return  $S_{BEST}$ 
8:     end if
9:   end for

```

725 The proposed **Algorithm 1** is used to pick the best simulated multiplexes based on the developed metrics
 726 ADS and MDS, and these assays are further tested empirically to select the optimal one for the diagnostic
 727 use. Moreover, to verify the correlation of the Smart-Plexer ranking with the ACA performance, the
 728 algorithm was used to select the bottom multiplexes with the lowest ADS and MDS, by modifying step 3 and
 729 4, so that the smallest instead of the largest ADS and MDS are applied.
 730
 731

732 **The Smart-Plexer Workflow.** The complete workflow of utilising the Smart-Plexer in a real laboratory setting
 733 is illustrated in **Figure 1** and depicted as follows: given a number of target genes to be identified, several
 734 candidate primer sets are first *in-silico* designed and tested in singleplex format for each target, resulting in
 735 real-time PCR amplification curves for all the assays. The obtained data are further processed using the
 736 background, late curve, and noisy curve removal techniques mentioned in the **Data Processing** section. The
 737 processed curves are then fitted with the sigmoidal function from which the “*c*” parameters are extracted.
 738 For each potential combination of primer sets, inter-target distances of “*c*” parameters from singleplex
 739 curves are calculated and function as simulated alternatives for empirical multiplex curve distances. In this
 740 way, the best candidates for multiplex assays can be selected by choosing the combinations with the most
 741 distant target clusters (represented by “*c*”) in the simulation. This progress is achieved by calculating the “*c*”
 742 parameter-based ADS and MDS of each combination and finding the best ones using the ranking system
 743 mentioned above. The best candidate assays shortlisted from simulated multiplexes further go through wet-
 744 lab tests on synthetic DNA templates, and the ACA-based target identification is applied to the empirical
 745 multiplex data. The final winner assay with the highest ACA classification performance on synthetic DNA is
 746 labelled as the optimal assay, which is the final output of the entire Smart-Plexer workflow.

747
 748 **3-plex validation.** Synthetic DNA of Adenovirus (HAdV), Human coronavirus HKU1 (HCoV-HKU1) and Middle
 749 East respiratory syndrome-related coronavirus (MERS-CoV) targets were selected for a 3-plex validation, and
 750 all the data were generated in real-time digital PCR (qdPCR). Three primer sets were designed as candidates
 751 for each target, resulting in 27 potential combinations of multiplex assays in total. Because of the relatively
 752 small number of candidate assays, it is possible to perform wet-lab experiments for all combinations and
 753 analyse the relationship between simulated and empirical multiplex curve distances. Simulated ADS and
 754 MDS were calculated on different levels of curve representations (raw curves, FFI-normalised curves, and
 755 fitted parameters), and their correlations with the same metrics derived from empirical multiplex data were
 756 analysed. Furthermore, the ADS and MDS of “*c*” parameters, which are more concise indicators for inter-
 757 target curve distances, were generated and compared between simulated and empirical multiplexes. ACA
 758 performance against simulated ADS and MDS was depicted, and the t-SNE of the selected assays’ results
 759 were illustrated.

760
 761 **7-plex validation.** Following the 3-plex validation, seven targets were used to further validate the Smart-
 762 Plexer performance, where each target had at least two different assays, resulting in a total of 24
 763 singleplexes and 4,608 candidate combinations. Unlike for 3-plex, the mass number of combinations makes
 764 it impossible to empirically test all the assays in multiplex settings. Instead, representative groups of assays

765 were chosen for the laboratory validation. Following the aforementioned Smart-Plexer workflow, after
 766 calculating simulated ADS and MDS on “c” parameters, six highest ranked (BEST) and six lowest ranked (BOT)
 767 combinations were picked out using the Ranking System. In addition, six middle-distant combinations (MID)
 768 were selected with the step below:
 769

Algorithm 2

- 1: **Initialise** N_{MID} as required, $S_{MID} \leftarrow \emptyset$, MDS_{max} and ADS_{max} the maximum MDS and ADS among all combinations, $ADS_{bias} = MDS_{bias} \leftarrow 0.001$
 - 2: $R_{MDS} \triangleq \left(\frac{MDS_{max}}{2} - MDS_{bias}, \frac{MDS_{max}}{2} + MDS_{bias} \right)$
 - 3: $R_{ADS} \triangleq \left(\frac{ADS_{max}}{2} - ADS_{bias}, \frac{ADS_{max}}{2} + ADS_{bias} \right)$
 - 4: $S_{MID}^{tmp} \triangleq \{x \mid MDS_x \in R_{MDS} \text{ and } ADS_x \in R_{ADS}, \forall x \in S_{ALL}\}$
 - 5: $S_{MID} \leftarrow$ apply **Algorithm 1** on S_{MID}^{tmp} with N_{MID}
 - 6: **return** S_{MID}
-

770 TOP-ADS and TOP-MDS (N=6) assays were selected empirically with large ADS but small MDS, and large MDS
 771 but small ADS, respectively. Similarly to the 3-plex validation, the relationship between simulated and
 772 empirical scores of the selected assays was explored by correlations of simulated and empirical metrics and
 773 comparisons of “c” parameter distributions. ACA was also applied to different groups of combinations. The
 774 complete pipeline of the 7-plex validation is illustrated in **Supplementary Fig. 3**.

775

776 **Clinical sample classification.** To verify the feasibility of Smart-Plexer in real clinical settings, we chose the
 777 optimal multiplex assay (PM7.2151) that achieved the highest ACA accuracy in synthetic DNA testing and
 778 conducted experiments on clinical samples. The multiplex was tested on the clinical samples using qdPCR,
 779 with 770 unprocessed raw amplification data (including flat curves) as the output for each sample. After the
 780 data processing step, the curves were input into an ACA classifier pre-trained with synthetic DNA data, and
 781 curve-level predictions were assigned to every positive curve. The target category of a sample was then
 782 decided by finding the mostly shown label among all the sample’s curve predictions. The confidence
 783 level of prediction is defined as the percentage of curves with this most shown label. Correctly predicted
 784 samples are marked as “detected”, otherwise “undetected”.

785 Acknowledgements

786 This work was supported by the Imperial COVID-19 Research Fund (WDAI.G28059); the Department of
 787 Health and Social Care-funded Centre for Antimicrobial Optimisation (CAMO) at Imperial College London;
 788 the Imperial College President’s PhD Scholarships 2021 (KX), the Imperial College’s Centre for Antimicrobial
 789 Optimisation (CAMO). Authors FB, KHC, AH, PG and JRM are affiliated with the NIHR Health Protection
 790 Research Unit (HPRU) in Healthcare Associated Infections and Antimicrobial Resistance at Imperial College
 791 London in partnership with the UK Health Security Agency (previously PHE) in collaboration with, Imperial
 792 Healthcare Partners, the University of Cambridge and the University of Warwick. AH is a National Institute
 793 for Health Research Senior Investigator.

794

795 References

- 796 44. QIAamp Viral RNA Mini Handbook - QIAGEN.
 797 [https://www.qiagen.com/gb/resources/resourcedetail?id=c80685c0-4103-49ea-aa72-](https://www.qiagen.com/gb/resources/resourcedetail?id=c80685c0-4103-49ea-aa72-8989420e3018&lang=en)
 798 [8989420e3018&lang=en](https://www.qiagen.com/gb/resources/resourcedetail?id=c80685c0-4103-49ea-aa72-8989420e3018&lang=en).
- 799 45. GenBank Overview. <https://www.ncbi.nlm.nih.gov/genbank/>.
- 800 46. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
 801 *Acids Research* **32**, 1792–1797 (2004).

802 47. Geneious | Bioinformatics Software for Sequence Data Analysis. *Geneious*
803 <https://www.geneious.com/>.

804 48. OligoAnalyzer Tool - primer analysis | IDT. *Integrated DNA Technologies*
805 <https://eu.idtdna.com/pages/tools/oligoanalyzer>.

806 49. Multiple Primer Analyzer - UK. [https://www.thermofisher.com/uk/en/home/brands/thermo-](https://www.thermofisher.com/uk/en/home/brands/thermo-scientific/molecular-biology/molecular-biology-learning-center/molecular-biology-resource-library/thermo-scientific-web-tools/multiple-primer-analyzer.html)
807 [scientific/molecular-biology/molecular-biology-learning-center/molecular-biology-resource-](https://www.thermofisher.com/uk/en/home/brands/thermo-scientific/molecular-biology/molecular-biology-learning-center/molecular-biology-resource-library/thermo-scientific-web-tools/multiple-primer-analyzer.html)
808 [library/thermo-scientific-web-tools/multiple-primer-analyzer.html](https://www.thermofisher.com/uk/en/home/brands/thermo-scientific/molecular-biology/molecular-biology-learning-center/molecular-biology-resource-library/thermo-scientific-web-tools/multiple-primer-analyzer.html).

809 **Author contributions**

810 L.M. and J.R.M. conceived, designed and supervised the project. L.M. and Y.C., developed the algorithm and
811 performed the data analysis. L.M., Z.L. and N.D. performed wet-lab experiments. K.X. review algorithm and
812 data analytics concepts. L.K., M.C.S., A.H. and P.G. contributed to the proofreading and editing of the
813 manuscript. L.M., K.X. and J.R.M. wrote the manuscript with input from all the authors.
814

815 Table 1 | Clinical validation results

Sample index	Panel ID (Randox, UK)	Expected Pathogen (True Label)	ACA Classified Pathogen (Predicted Label)	AC count	Confidence Level (%)	Outcome
1	QAV164189	HAdV	HAdV	14	100.0	detected
2	QAV164189	HCoV-NL63	HCoV-NL63	770	100.0	detected
3	QAV164189	HCoV-NL63	HCoV-NL63	545	96.15	detected
4	QAV164189	HCoV-OC43	HCoV-OC43	94	78.72	detected
5	SCV2QC	SARS-COV-2	SARS-COV-2	769	69.96	detected
6	SCV2QC	SARS-COV-2	SARS-COV-2	631	94.77	detected
7	SCV2QC	SARS-COV-2	SARS-COV-2	766	100.0	detected
8	SCV2QC	SARS-COV-2	SARS-COV-2	756	99.34	detected
9	SCV2QC	SARS-COV-2	SARS-COV-2	748	99.20	detected
10	QAV154181	MERS	MERS	287	60.98	detected
11	QAV154181	MERS	MERS	770	96.49	detected
12	QAV154181	MERS	MERS	770	79.09	detected
13	QAV154181	MERS	MERS	698	91.69	detected
14	QAV154181	MERS	MERS	20	70.00	detected

816
817 Table 2 | five sigmoidal fitted parameters

Parameter	Mathematical meaning	Representation in amplification curves
<i>a</i>	Amplitude of the function in the y-axis	Affect the maximum fluorescence that the amplification curves can reach
<i>b</i>	Vertical shift of the function along the y-axis	Affect the maximum fluorescence together with parameter <i>a</i>
<i>c</i>	Maximum slope of the sigmoid function	Related to the efficiency of PCR reactions
<i>d</i>	Horizontal shift of the function x-axis	Fractional cycle of the inflection point (related to C_t values)
<i>e</i>	Richard's coefficient	Asymmetry of the sigmoidal trend

818

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementarydata.docx](#)