

Citation for published version:

Huang, J, Li, J, Zhou, J, Wang, L, Yang, S, Hurst, L, Li, W-H & Tian, D 2018, 'Identifying a large number of high-yield genes in rice by pedigree analysis, whole genome sequencing and CRISPR-Cas9 gene knockout: High-yield genes detected by pedigree analysis', Proceedings of the National Academy of Sciences of the United States of America (PNAS). <https://doi.org/10.1073/pnas.1806110115>

DOI:

[10.1073/pnas.1806110115](https://doi.org/10.1073/pnas.1806110115)

Publication date:

2018

Document Version

Peer reviewed version

[Link to publication](#)

Copyright © 2018 (National Academy of Sciences). The final publication is available at Proceedings of the National Academy of Sciences via <https://doi.org/10.1073/pnas.1806110115>.

University of Bath

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Identifying a large number of high-yield genes in rice by pedigree analysis, whole genome sequencing and CRISPR-Cas9 gene knockout

Ju Huang¹, Jing Li¹, Jun Zhou², Long Wang¹, Sihai Yang¹, Laurence Hurst³, Wen-Hsiung Li⁴, Dacheng Tian¹

¹Nanjing University, ²Harvard University, ³University of Bath, ⁴Academia Sinica

Submitted to Proceedings of the National Academy of Sciences of the United States of America

Repeated artificial selection of a complex trait facilitates the identification of genes underlying the trait, especially if multiple selected descendant lines are available. Here we developed a pedigree-based approach to identify genes underlying the Green Revolution (GR) phenotype. From a pedigree analysis, we selected 30 cultivars including the "Miracle rice" IR8, a GR landmark, its ancestors and descendants, and also other related cultivars for identifying high-yield genes. Through sequencing of these genomes, we identified 28 ancestral chromosomal blocks that were maintained in all of the high-yield cultivars under study. In these blocks, we identified 6 genes of known function, including the GR gene *sd1*, and 123 loci with genes of unknown function. We randomly selected 57 genes from the 123 loci to do knockout or knockdown studies and found that a high proportion of these genes are essential or have phenotypic effects related to rice production. Notably, knockout lines have significant changes in plant height ($p < 0.003$), a key GR trait, compared to wild-type lines. Some gene knockouts or knockdowns were especially interesting. For example, knockout of *Os10g0555100*, a putative glucosyltransferase gene, showed both reduced growth and altered panicle architecture. In addition, we found that in some retained chromosome blocks, several GR related genes were clustered, although they have unrelated sequences, suggesting clustering of genes with similar functions. In conclusion, we have identified many high-yield genes in rice. Our method provides a powerful means to identify genes associated with a specific trait.

high yield gene | pedigree analysis | green revolution | gene knockout

Complex traits, which might be related to survival in natural environments or crop productivity (1), are genetically difficult to dissect. This is, in part, because the effect of a single gene on a phenotype is usually small (2). To determine the genetic architecture of a complex trait (and the underlying gene networks), the most commonly employed methods are quantitative trait loci (QTL) mapping and genome-wide association studies (GWAS). QTL mapping is suitable for relatively simple quantitative traits (3), while GWAS provides valuable insights into trait architecture or candidate loci (4). Both methods have limitations, however. In QTL, the effects detected may be sensitive to external environments (5) and the span of chromosomal regions detected is often too long (owing to limited recombination events (6)) to pinpoint the causative gene(s). Similarly, in GWAS, the effects detected are sensitive to population structure, leading to both false positives and false negatives (7, 8).

Recently, a pedigree from crosses between different founding genotypes was used to fine-map QTLs in *Arabidopsis* (1, 9). The pedigree-based analysis combines linkage and association study (6). A pedigree with a founding genotype (e.g., derived from a single cross of two ancestors) and with recombination events over many generations could overcome the disadvantages inherent in QTL and GWAS. To reduce the sensitivity to environmental effects, however, it is necessary to have a clear phenotypic difference between the two ancestors. Identification of chromosomal blocks preserved in all members of the pedigree under

selection for a given trait will facilitate identification of candidate genes. The question then is whether these candidates are indeed associated with the trait. The CRISPR-cas9 system (10) can in principle be used to knock out each candidate gene to get an insight into its function. Below we describe an application of this pedigree/knockout approach to the identification of high-yield genes in rice.

Our study takes advantage of the diploid rice pedigree in the Green Revolution. The Green Revolution has dramatically increased agriculture production worldwide since the 1960s, saving millions of lives from food shortage (11). The novel technologies allowed agronomists to breed high-yield varieties of maize, wheat, and rice. The yields were more than doubled in developing countries from 1961–1985 (12). Perhaps the most significant milestone of the Green Revolution was the introduction of semi-dwarfing genes into selected rice cultivars by hybridization.

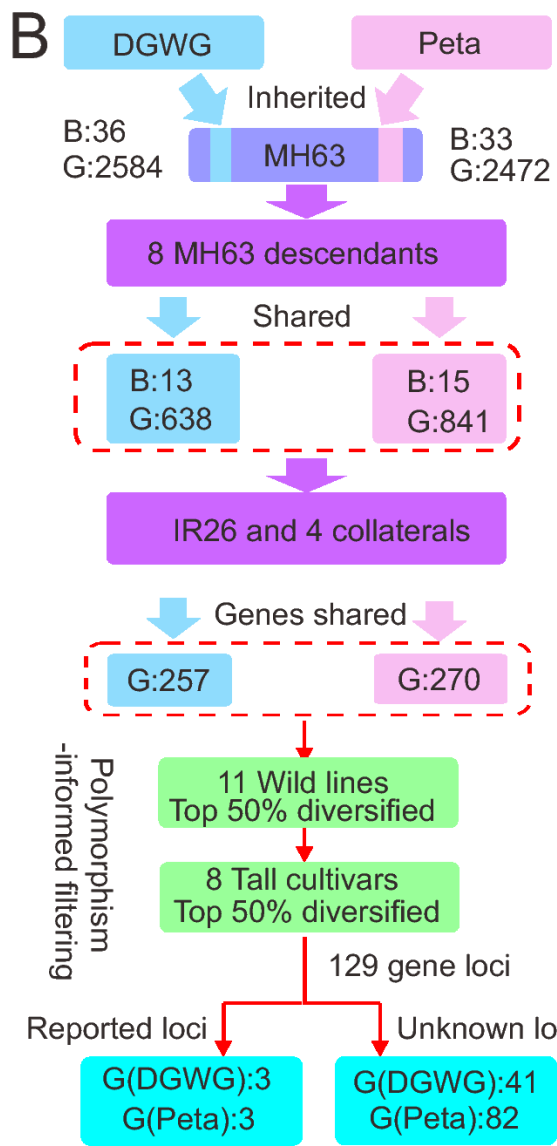
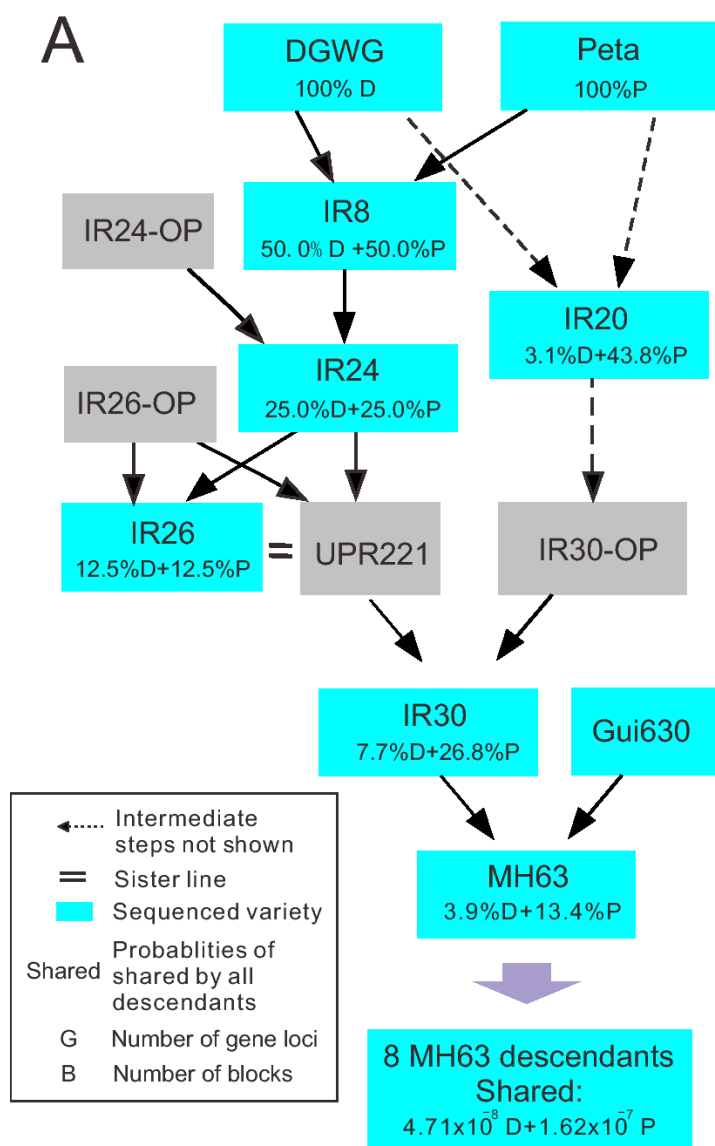
The first semi-dwarf and high-yield modern rice variety (HYV) of the Green Revolution, known as the "Miracle rice" IR8, was created by crossing the Indonesian variety "Peta" with the Chinese variety "Dee-geo-woo-gen" (DGWG). It represented the first generation of the "high-yielding plant type", which provided a significantly higher yield potential for irrigated rice (13). In addition to the significant reduction in stem length, the high-yield rice cultivars have other important traits such as an early flowering time, improvement in photosynthetic allocation, and insensitivity to day length, directly or indirectly influencing the grain yield and yield stability (14, 15). These high-yield traits could

Significance

Finding the genes that control a complex trait is difficult because each gene may have only minor phenotypic effects. Quantitative trait loci mapping and genome-wide association study techniques have been developed for this purpose but are laborious and time-consuming. Here we developed a new method combining pedigree analysis, whole genome sequencing and CRISPR-Cas9 technology. By sequencing the parents and descendants of IR8, the "Miracle rice" in Green Evolution, we determined many genes that had been retained in the pedigree by selection for high yield. Knockout and knockdown studies showed that a large proportion of the identified genes are essential or have phenotypic effects related to production. Our approach provides a powerful means for identifying genes involved in a complex trait.

Reserved for Publication Footnotes

137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204



205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272

Fig. 1. Pedigree and flowchart for the identification of geneloci under selection. (a) An abridged pedigree of the major rice cultivars used in this study. The green-highlighted cultivars were re-sequenced, while the gray-highlighted were not. OP means “the other parent” and was not sequenced. The percentage in a box shows the expected probability of a given locus inherited from DGWG (D) or Peta (P) in that generation. The bottom box indicates the expected probabilities of a locus shared by all of the 8 MH63 descendants, which are extremely low (see *SI Appendix*, Table S4). A solid arrow denotes a direct parent (i.e. IR20) and a dotted arrow indicates an indirect ancestor (i.e. IR24). (b) Flowchart of the approach used to identify candidate blocks and gene loci derived from DGWG or Peta. Numbers of blocks (B) and gene loci (G) within the high confidence blocks are shown in each step of filtering. The reported 6 genes (3 from DGWG and 3 from Peta) are the gene loci that have clear functions reported in literature. Most of the 129 gene loci each contain only one gene, except that 28 of them have two or more overlapped genes within a locus (see Methods).

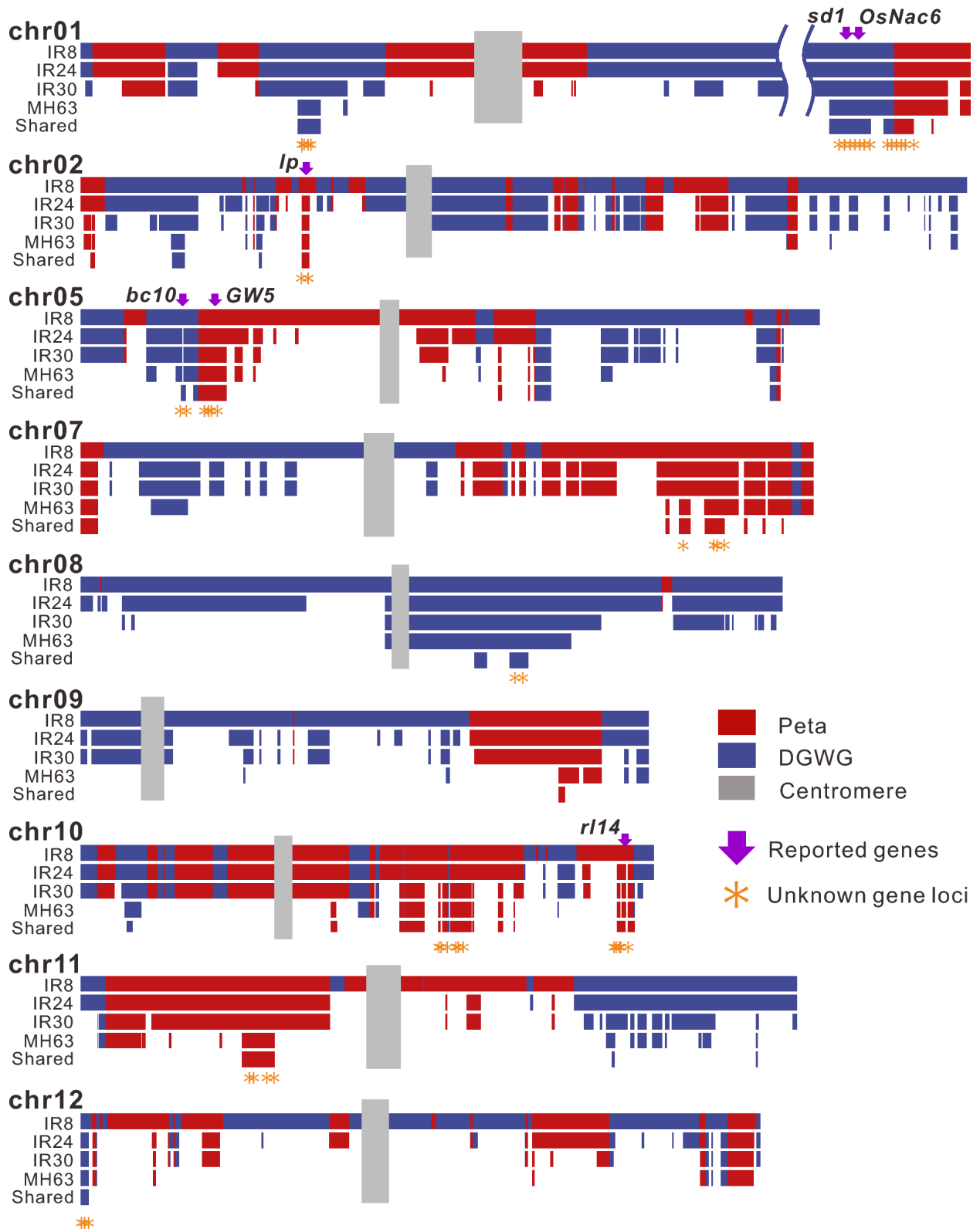
be traced from the pedigree of “Miracle rice” IR8 that consists of its parents and high-yield progenies.

We assume that the genes related to high-yield were under strong artificial selection because yield was the major target trait of rice breeding since the 1960s. In this scenario, we note: 1) if the multiple lineages descended from an original cross have all been placed under the same selection, the alleles responsible for the trait in question should be found in all of the descendants, but not in all control populations; 2) in principle these alleles can be traced back to their origination and any variants inherited in all generations can be identified; 3) a gene under strong artificial selection should be present more commonly in progeny than genes not under selection; and 4) when knocking out a high-yield gene, a changed plant phenotype (e.g., an observable change in morphology or physiological response such as sterility) should

be observed. All these expectations can be tested by sequencing the cultivars at important nodes in the pedigree and then by a knockout study using the CRISPR-Cas9 system.

Using the strategy above, we studied the extended pedigree of the ancestors and descendants of IR8 and other related lines (Fig. 1A) to determine a set of genes that played a critical role in the rice Green Revolution. By resequencing 30 cultivars from the pedigree (Fig. 1), we identified 28 chromosomal blocks, including 129 candidate gene loci, that have been preserved by artificial selection (Fig. 2). Fifty-seven gene loci with unknown function were selected to do knockout using the CRISPR-Cas9 technique. If the knockout failed, then a knockdown experiment was conducted. We found that 79% (15/19) knocked out loci and 62% (10/16) knocked down loci have phenotypic changes. These studies revealed a striking enrichment in yield/morphology-associated

273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340



341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408

Fig. 2. Blocks inherited from DGWG and Peta in IR8, IR24, IR30, MH63 and the 8 descendants of MH63. Blue and red bars represent blocks derived from DGWG and Peta, respectively. "Shared" denotes the shared regions in all of the eight MH63 descendants. The purple arrows represent the 6 genes reported with functions related to plant type or high-yield, while the asterisks represent the 123 gene loci with unknown functions; the 6 genes are shared by all 8 MH63 descendants and 5 collateral series. Chromosomes 3, 4 and 6, which contain no regions shared by all 8 MH63 descendants, are not shown here. The second last block on chromosome 1 was shortened using breaks.

genes among the candidate genes. Thus, our pedigree-guided

approach provides a simple, robust and fast means to identify candidate genes under directional selection.

Table 1. Numbers of blocks derived from DGWG and Peta in different descendants

Ancestor	Descendant	Chromosome												Total
		1	2	3	4	5	6	7	8	9	10	11	12	
DGWG	IR8	4	15	17	12	7	11	4	3	3	13	6	13	108
	IR24	4	12	9	8	6	11	4	3	3	13	3	9	85
	IR30	4	11	7	7	6	7	4	2	3	12	2	8	73
	MH63	2	5	3	2	4	2	2	1	3	7	2	3	36
	Shared ^a	2	2	0	0	3	0	0	1	0	3	1	1	13
Peta	Genes ^b	442	75	0	0	136	0	0	64	0	28	6	34	785
	IR8	4	15	17	12	6	11	5	2	2	12	5	12	103
	IR24	4	14	11	6	4	6	5	1	2	10	3	10	76
	IR30	4	12	10	6	4	2	5	0	2	10	3	8	66
	MH63	1	5	7	1	3	0	3	0	1	6	1	5	33
	Shared ^a	1	2	0	0	3	0	2	0	1	5	1	0	15
	Genes ^b	101	34	0	0	115	0	265	0	42	308	95	0	960

^a Shared represents the blocks and enclosed genes observed in MH63 and in all of its eight descendants.

^b Genes contained in the shared blocks.

Submission PDF

Table 2. Phenotype when a specific gene was knocked out

Sampled ancestral block	Loci	Observed phenotypes
DGWG chr01:37602014-39226171	Os01g0884200	Dwarf, sterile
	Os01g0884400 ^a	Late heading, sterile
	Os01g0884450	
	Os01g0885000	Small, growth retarded, fewer tillers
Peta chr01:40248759-40971796	Os01g0886000	Late heading, fewer tillers, sterile
	Os01g0925600 ^a Os01g0925700	Rolling leaves, shorter panicle, dwarf
	Os01g0930800	Late heading, sterile
	Os01g0930900	No phenotypic change
	Os10g0555600 ^a	Dwarf
Peta chr10:21769689-21922126	Os10g0555651	
	Os10g0555900 ^a Os10g0556000	Dwarf, late heading
	Os10g0556200	Dwarf
	Os10g0556900	No phenotypic change
	Os10g0555100	Dwarf, spike shape change ,
	Os10g0555200	Dwarf, sterile
	Os10g0555300	Dwarf, sterile
	Os10g0555700	Sterile
	Os10g0556100	Small, growth retarded, leaf rolling
	Os10g0558850	Rolling leaves, dwarf, weak
Os10g0559800 ^a Os10g0559833	No phenotypic change	
Peta chr10:21992900-22072751		
Peta chr11:6540176-7824094	Os11g0242400	No phenotypic change

The 123 gene loci that passed our filtration came from 16 blocks, which ranged in size from 43kb to 1624kb. In total, 19 gene loci from 5 blocks of different sizes (80kb-1624kb) were successfully knocked out. For each gene, about 15 independently transgenic plants were obtained and on average in 79.5% of the cases the gene was knocked out in both homologous chromosomes. The phenotypic change was based on the observation of the homozygous knockout plants. No phenotypic change means no significant change in phenotype; e.g., the knockout of the locus Os01g0930900 showed shorter plants and shorter awns, but the changes were not statistically significant. In total, 15 out of the 19 knockouts exhibited phenotypes different from the wild type, suggesting that a large portion of these unknown-function gene loci are involved in flowering, fertility, leaf morphology, etc. The genotype and phenotype of each gene studied are described in Table S15. All the knockout plants in this table were in the Kasalath background

^aIn five pairs, the two genes in a pair are partly or completely overlapped. For example, Os01g0884450 is completely contained in Os01g0884400.

Table 3. Phenotypic changes in knock-down mutants.

Locus	Abnormal phenotypes
Os01g0883900	Curled leaves, retarded growth. Died before matured.
Os01g0931600	Retarded growth, multiple tillers.
Os05g0170200	Retarded growth, curled leaves.
Os10g0556500	Brown and curled leaves. Died before matured.
Os10g0556700	Normal.
Os10g0559866	Normal.
Os02g0258900	Retarded growth, brown and curled leaves.
Os10g0391100	Normal.
Os10g0391200	Normal.
Os10g0392400	Curled leaves. Died before matured.
Os10g0554900	Normal.
Os12g0103000	Brown leaves. Died before matured.
Os12g0104250	Normal.
Os12g0104400	Brown leaves. Died before matured.
Os12g0104700	Retarded growth, curled leaves. Died before matured.
Os12g0104733 ^a	Only grew roots. No seedling.
Os12g0104766 ^a	Curled leaves.

^a These genes are included in the same locus.

Results

Rice cultivar selection and SNP identification

The famous “Miracle rice” IR8 is the key cultivar in our pedigree analysis (Fig. 1A). Its descendants and derivatives have been extensively used in the field, and its parents have been widely utilized to breed desired plant types (16). Another key cultivar is Minghui63 (MH63), which is a fourth generation descendant of IR8 and was the restorer line for a number of rice hybrids. MH63 accounted for >20% of the total production area in China during the 1980s and 1990s (17). Because of its wide planting areas with a stable high-yield, environmental or epigenetic effects could be excluded. IR8 and MH63 form the basis of our pedigree analysis. The pedigree further expands upward to the parents of IR8 (i.e., DGWG and Peta) and MH63 (IR30 and Gui630) and downward to the descendants of IR8 (i.e., IR24) and MH63. IR20, which has the same parents as IR8, and eight extensively used descendants of MH63 are also included in the analysis (Fig. 1A). All descendants of IR8 possessed the common feature of high-yield. To enhance the resolution in identifying genes under selection, we also sequenced four IR8 collateral series, eight tall landraces and a wild rice as the controls (SI Appendix, Fig. S1B and Table S1). The alleles present in the control groups were considered unlikely to contribute to high yield.

The 30 diploid rice accessions selected above were re-sequenced with a reasonable coverage depth (>20×) in our study (SI Appendix, Table S1 and Fig. S2). Because pedigree information and independent re-sequencing of descendants from the same ancestor offer the unique advantage of discriminating against false markers, each inherited block of interest can be double-checked not only between successive generations but also between nodes separately by more than one generation and between lineages. Based on the linked markers in the majority of the successive generations, this approach can exclude false markers, infer correct single nucleotide polymorphisms (SNPs), and improve the accuracy of SNP identification (SI Appendix, Fig. S3). In the two most important parent-offspring trios, DGWG-Peta-IR8 and IR30-Gui630-MH63, a total of 592,603 and 481,385 high-quality SNPs were called, respectively, to detect the inherited chromosomal blocks from IR8 and its parents (see Methods and SI Appendix, Fig. S4).

Expected and observed proportions of inherited blocks

With the pedigree information, the probability of a block or a gene being passed on to the next generation can be computed using classical genetic theory. One can then compare the computed probability with the observed proportion (see SI Appendix, SI Materials and Methods). In the absence of selection, the probabilities of a gene locus in MH63 from DGWG and Peta are expected to be 3.9% and 13.4%, respectively (Fig. 1A and SI Appendix, Table S2). The probability of one or more DGWG or Peta blocks being present in all eight descendants of MH63 is extremely low (4.71×10^{-8} or 1.62×10^{-7}) (Fig. 1A, Table 1 and SI Appendix, Table S3, Table S4). Therefore, every block retained in all of the MH63 progenies is likely to have been targeted by artificial selection for the high-yield phenotype.

Theoretically, the heterozygosity of the F₁ hybrid will be reduced to half in its F₂ progeny through selfing and will eventually be reduced to almost zero in an inbred line (e.g., IR8 or MH63). Therefore, the crossover events can be detected in both IR8 and MH63 to determine the origin of each block (SI Appendix, Table S5 and Table S6). The block information in MH63 enabled us to exclude the genetic blocks from Gui630 and identify those from DGWG or Peta based on the pedigree in Fig. 1A. In MH63, we found 57 and 59 blocks that were derived from 36 DGWG and 33 Peta blocks in IR8, respectively (Fig. 2). Thus, many of the original inherited blocks from DGWG and Peta had been fragmented into smaller ones in MH63 by recombination. The average length is 483 kb for the 57 DGWG-derived blocks and 398 kb for the 59 Peta-derived blocks, which are 5.45- and 3.20-fold shorter than the average lengths in the original blocks in IR8, respectively (SI Appendix, Table S7 and Table S8). Among those original blocks, only a total of 6.26 Mb DGWG and a total of 8.76 Mb Peta segments are inherited in all of the 8 MH63 descendants. They were 2.39- and 1.55-fold shorter than the inherited blocks observed in MH63, respectively. The sequences shared by all 8 MH63 descendants contained 785 DGWG- and 960 Peta-specific genes (Fig. 1B and Fig. 2).

Identification of candidate genes for the high-yield phenotype

When only a limited number of genes in a block are under selection, the ancestral block will become shorter and shorter over generations because of recombination events. Fig. 2 includes an example in which a block on Peta chromosome 5 became shorter and shorter by crossover events from IR8 to MH63. Interestingly, a candidate gene, *GW5*, which is responsible for rice grain width, shape, quality and yield, is located near recombination hotspots (18) but has been retained. The pattern displays efficient selection on this block.

In a block with many genes, some alleles that are not subjected to selection may be inherited due to linkage (i.e., hitchhiking). Several strategies were employed to exclude the hitchhiked genes and identify the genes that were most likely the target of selection, including those with un-annotated functions (Fig. 1B). π (polymorphic sites/informative sites) was calculated for each 10-kb window to compare the diversity values within and between different groups. First, we assumed that targeted alleles should have been retained also in the IR8 collateral series because those cultivars are also of high-yield plant types. With this assumption, we selected four cultivars of the IR8 collateral series (SI Appendix, Fig. S1 and Table S1) and calculated the nucleotide diversity of these candidate genes between MH63 and each of the four collateral cultivars together with IR26, a progeny of IR24, and a sister line of UPR221 (a parent of IR30 in Fig. 1A). Only the genes that had an average diversity <0.0001 and were identical in the majority of collateral series (≥ 3) for the compared pairs were kept. Second, we assumed that a gene with an extremely low diversity among wild rice lines and cultivars should be excluded because it is more likely to be essential for fundamental biological processes rather than being responsible for the high-

681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748

749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816

Table 4. Plant height comparison between Kasalath knockout mutants and wild types.

	Locus	Mutant height (cm)	WT height (cm)	Height change: (mutant height-WT height)/WT height
4 positive controls ^a	Os01g0883800	62.8	132.4	-52.6%
	Os01g0884300	65.5	129.3	-49.3%
	Os05g0170000	47.8	130.3	-63.4%
	Os02g0260200	98.3	123.3	-20.3%
18 target gene loci ^b	Os01g0884200	110.7	129.3	-14.40%
	Os01g0884400	125.3	127.3	-1.6%
	Os01g0884450			
	Os01g0886000	127.7	129.6	-1.5%
	Os01g0925600	93.3	125.6	-25.7%
	Os01g0925700			
	Os01g0930800	125.3	131.2	-4.5%
	Os01g0930900	130.2	129.1	0.9%
	Os10g0555100	99.6	130.1	-23.4%
	Os10g0555200	99.7	130.3	-23.5%
	Os10g0555300	105.2	129.2	-18.6%
	Os10g0555600	95.2	130.2	-26.9%
	Os10g0555651			
	Os10g0555700	127.1	130.6	-2.7%
	Os10g0555900	64.6	132.4	-51.2%
	Os10g0556000			
	Os10g0556100	65.2	127.3	-48.8%
	Os10g0556200	73.7	122.1	-39.6%
Os10g0556900	131.2	129.1	1.6%	
Os10g0558850	117.2	130.2	-10.0%	
Os10g0559800	130.1	129.3	0.6%	
Os10g0559833				
Os11g0242400	129.4	128.6	0.6%	
10 random controls ^c	Os01g0936100	130.3	131.5	-0.9%
	Os05g0375600	134.0	132.4	1.2%
	Os05g0571700	126.5	129.2	-2.1%
	Os05g0573600	132.0	130.1	1.5%
	Os10g0341750	134.0	130.1	3.0%
	Os10g0342300	132.0	129.2	2.2%
	Os10g0341700	133.0	130.2	2.2%
	Os05g0571300	134.3	132.4	1.5%
	Os10g0558400	128.5	132.4	-2.9%
	Os10g0342650	131.3	131.5	-0.1%

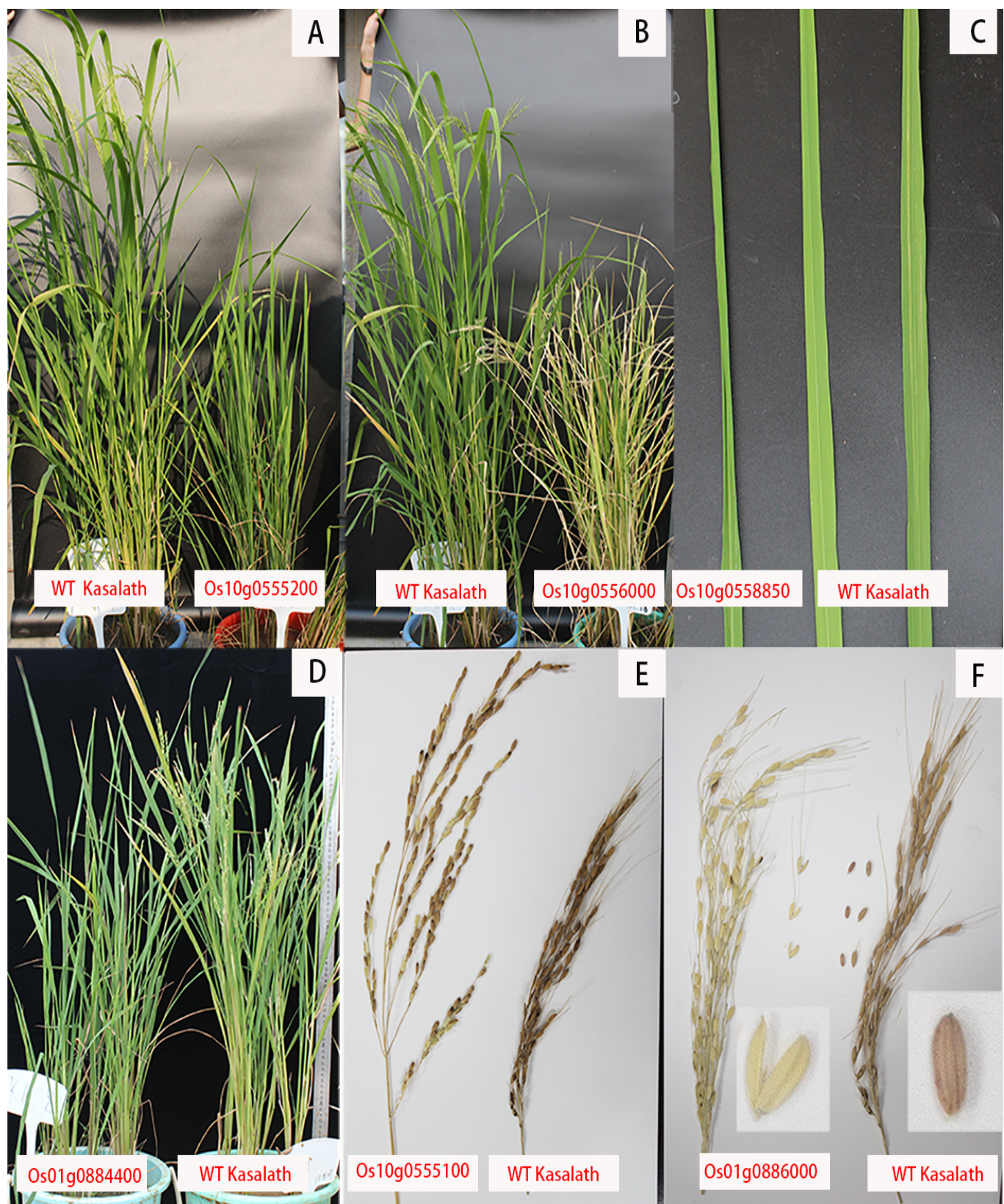
On average, positive controls showed 46.4% reduction in plant height ($p=0.017$, two-tail t-test, 95% confidence interval: -99.6 -20.84), while 18 target gene loci showed 16.4% reduction in plant height ($p=0.0013$, two-tail t-test, 95% confidence interval: -31.98 -9.22). 10 random controls only showed a slight difference (knockout effect) (average 0.5%, $p=0.42$, two-tail t-test, 95% confidence interval: -1.16 - 2.54).

^a Four of the 6 positive controls were knocked out in Kasalath. GW5 was knocked out in Wuyungeng and r114 was not successfully knocked out.

^b The knockout plant (Os01g0885000) died before the tillering stage, and the plant height could not be compared with the others. Therefore, only 18 target mutants were measured.

^c 10 genes adjacent to the target blocks were randomly chosen as controls.

817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884



885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952

Fig. 3. Photos of knockout mutants with changed phenotypes. These photos show 6 examples as shorter plants, rolling leaves, a later heading date, changed panicles and empty seeds compared with the wild type. The other 9 knockout mutants with observable phenotypic changes and the controls are shown in *SI Appendix, Fig. S7*. **Supporting information** The following materials are available in the online version of this article. **SI Materials and Methods**

yield phenotype. Therefore, we further filtered out the bottom 50% of genes in terms of the diversity between MH63 and the

11 wild rice varieties. Third, we filtered further by comparison to tall cultivars as follows. All the re-sequenced cultivars in this study were grown in the field and their heights were measured. Because the semi-dwarfism trait was specifically selected for the Green Revolution, we expect that the alleles related to the Green Evolution would be divergent with tall cultivars and only be kept in the genes showing a diversity higher than the median between MH63 and each of the 8 tall cultivars (*SI Appendix*, Table S1 and Table S9).

The above filtering procedure identified 129 gene loci, which can be divided into 101 single loci and 28 loci with overlapping genes (where two or more genes overlap completely or partly within the same locus). As an example of overlapping genes, the coding sequence of Os01g0883850 is completely contained in the reported gene *sd1* (Os01g0883800). These two genes are thus considered as a single entity in our analysis. Each locus is named by one gene it contains. Of the 129 gene loci, 44 are from DGWG- and 85 from Peta-specific blocks (Fig. 1B and *SI Appendix*, Table S10). These 129 gene loci are located on 17 blocks which are inherited in all 8 descendants of MH63. Six of the 129 gene loci contain genes with known functions, including the semi-dwarf gene, *sd-1*, known as the “green revolution gene”. This gene encodes gibberellin 20-oxidase, the key enzyme in the gibberellin biosynthesis pathway. Another gene, *larger panicle* (*lp*), which controls the panicle architecture (19), has recently been found to be a target of selection in Indica cultivars by a GWAS study of 1479 rice accessions (20). The others are *GW5*, *bc10*, *r14* and *OsNAC6*, responsible for grain width, brittle culm, leaf rolling and stress tolerance, respectively (18, 21–24). Interestingly, half of these six genes were identified from natural mutants in contrast to the fact that most functional genes were identified commonly from T-DNA insertion and mutagen induced mutants (roughly accounting for 90% of genes reported with known function). This suggests that the identified genes from a pedigree analysis could better reflect the real targets of selection in plant breeding than the genes identified from artificial mutants.

Knockout phenotypes of candidate gene loci

To determine whether a gene locus with unknown function has a phenotypic effect when knocked out, 57 of the 123 loci with unknown function were randomly sampled to do knockout by the CRISPR-cas9 system. Of these, 19 had knockout mutants, which were confirmed by PCR and Sanger sequencing. However, in the other 38 gene loci, no knockout mutants were obtained even after at least two independent transformations. We suspected that many of these genes are essential in callus development, so that no transformant survived. This possibility is supported by the observation that most (91.2%) of these genes had medium or high expression levels in callus (*SI Appendix*, Table S11).

As positive controls, we also attempted to knock out the 6 genes with known functions. As expected, 5 knockout mutants exhibited similar or stronger phenotypic changes compared to previous studies (18, 19, 21, 23, 24) (*SI Appendix*, Table S12). However, one of them, RL14, had no knockout mutant (see *SI Appendix*, Table S12). In a previous report, *r14*, which carries a single amino acid mutation, exhibited severe leaf rolling and therefore RL14 may have essential functions, so that its knockout could not survive (22). In addition, as random controls, 10 genes were randomly sampled from the 1kb-300kb regions (*SI Appendix*, Table S13) adjacent to the retained ancestor blocks (which were shared by all 8 descendants of MH63). The near-neighbor controls may be considered as conservative random controls for, unlike true random controls, these controls in part allow for possibly important position effects, such as the clustering of genes with similar expression profiles (25). In all 10 cases the knockout mutant showed no phenotypic changes (*SI Appendix*, Table S14), in contrast to 79% (15/19) of the unknown gene loci that showed observable phenotypic changes when the gene was knocked out

(Table 2 and detailed changes in phenotypes and genotypes in *SI Appendix*, Table S15).

High yield plants are typically dwarf, as dwarfism reduces investment into stalk, thereby potentially increasing investment into seeds. Therefore, we studied the growth difference between the mutated and unmutated version. We compared plant heights in knockout and wild type lines by the paired t-test (Table 4). As expected, the random control genes showed no difference in height between mutant and non-mutant versions ($P = 0.42$, 95% confidence interval, -1.15cm, 2.54cm), while the positive controls showed a significant shorter height in mutants than in wild type ($P = 0.017$, confidence interval: -99.6 cm to -20.84 cm). Importantly, for the test group we also saw a strong dwarfism phenotype ($P = 0.0013$; 95% CI: -31.98 to -9.22 cm). As these were random samples from the 123 unknown gene loci, it implies that a high proportion of the 123 loci have a phenotype similar to that of the well-described positive control genes identified by the same method. However, the extent of the dwarfism is reduced in the test sample compared with the positive controls (*t*-test on percentage difference comparing positive control and test samples, $P = 0.029$, 95% CI: -67.82 -5.48). These genes may have weaker effects than the previously reported ones, and this may be why they have not been identified.

A gene of particular interest is Os10g0555100, as its knockout showed a different panicle architecture and a 23% reduction in height. Note that one of the reported genes, *larger panicle* (*lp*), showed an altered panicle architecture as well. The protein product may be a glycogenin glucosyltransferase (see ic4r.org), suggesting a possible role in controlling free glucose and glucose storage. This speculation, however, requires further analysis. Among the other genes some, such as Os10g0558850, had rolled leaves (Fig. 3) but a relatively modest reduction in plant height (~ 10%). All the 15 unknown gene loci with knockout phenotypes have various protein-level motifs with unknown function, suggesting that the plant type and the high-yield phenotype are controlled by many types of genes.

Interestingly, the physically-proximal gene loci, although showing no sequence similarity, have similar functions. For example, 3 of the 6 gene loci on chromosome 1 (from Os01g0884400 to Os01g0930900) had knockouts resulted in late heading and 6 of the 11 loci on chromosome 10 (Table 2 and *SI Appendix*, Table S15) had knockouts resulted in dwarf phenotypes relative to the background line. This clustering mirrors the previously observed clustering of QTL signals (26). The clustering may reflect selection for coordinated gene expression or may possibly be owing to epistatic effects. Importantly, this result also suggests a strategy for finding genes with similar functions: if you have found one, investigate its neighbors.

Knockdown phenotypes of gene loci with no knockout transformant

To investigate the 38 loci with no knockout mutants, we randomly selected 26 loci to knock down their expression level, using the dCas9 knockdown technique (27). Similar to the knockout results, in 10 of the 26 loci (38.5%) no knockdown mutants were obtained due to the death of the transformed callus after hydromycin selection. Most of the 26 loci also have medium or higher expression levels in callus (*SI Appendix*, Table S16). Moreover, even in the 16 loci with knockdown transgenic plants, 10 knockdown plants showed distinct negative phenotypic changes and 7 of them died during plant regeneration (Table 3 and *SI Appendix*, Fig. S5). As expected, the qRT-PCR study confirmed that the expression of these target loci in knockdown transformants was indeed down-regulated (*SI Appendix*, Fig. S6). These results suggest that most of the 38 candidate genes are essential genes in rice.

Discussion

Determining the genes that explain complex traits has never been easy. The two much used methods, QTL and GWAS, have both led to important discoveries, but such analyses are typically very labor intensive. Indeed, during the past decades, much effort has gone into dissecting the genetic basis of high-yielding traits based on molecular linkage maps, e.g., the identification of many quantitative trait loci (QTLs) (28–31), but had identified relatively few genes. The pedigree-based method that we expanded here has, for some cases, short-cut much of the effort. It requires a good pedigree and consistent directional selection, however. Confirmation of such results would until recently also have been very time consuming, but CRISPR-Cas9 can greatly reduce the amount of work. In this study. We have not only identified the 3 well known loci for the Green Revolution (the “green revolution gene” *sd1*, grain size-related gene *GW5* and domestication gene *lp*), but also identified over 100 candidates. Among the 57 candidate genes selected for knockout and knockdown studies, we found that many of them are essential genes or showed phenotypic effects. Thus, the pedigree approach seems to be highly efficient for identifying candidate genes that were subject to strong selection.

While the knockout analysis suggested a low false positive rate, the false negative rate, by contrast, is unknown and probably quite high as our filters are quite stringent. Indeed, when we look at two genes that failed to pass the diversity cutoff, we find that one of them resulted in phenotypic change when knocked out. This suggests that slight relaxation of the stringent filtering will result in more candidates, but potentially a higher false positive rate too. More generally, we do not know how many genes are essential for the rice Green Revolution. As a consequence, the method should then be considered a technique to greatly enrich for selectively relevant genes rather than a method for an exhaustive search.

This study showed that rice is unusually well suited to this pedigree method. First, the well-documented pedigree information can be used to calculate the expected proportions of blocks (or loci) being transmitted from an ancestor to a descendant (e.g., Fig. 1A). By comparing the expected and observed proportions, the gene loci that were most likely to have been the target of artificial selection could then be identified. For example, the probability of a DGWG block appearing in all of the eight MH63 descendants was estimated to be nearly zero. Thus, if a block is observed in the re-sequencing data, it was very likely subjected to strong artificial selection. Second, from the relationships in a pedigree, SNP markers can be verified and corrected by comparing the sequences of parents and offspring between generations (demonstrated in Fig. 2 and *SI Appendix*, Fig. S1). In rice we are fortunate to have access to the stocks of the prior generations. Third, pedigree analysis focuses on tracing relatively longer blocks from the parents to the offspring instead of single SNPs or genes. It is therefore not difficult to identify selected targets. Finally, the CRISPR-Cas9 system provides an effective way of gene knockout to find a set of genes relevant to complex traits. In conclusion, our approach should be useful for many breeding projects.

Our choice of our model organism was not just motivated by the fact that the conditions for pedigree analysis were met, but also by the enormous impact of the Green Revolution, as indicated by the generation of high-yielding plant types through breeding. The introduction of dwarfing genes has resulted in plants that possess short and strong stalks, which are less liable to lodging. The stability of shorter plants dramatically reduces the need for photosynthetic investment in the stem. Assimilates are then redirected to grain production, resulting in a better plant type and increased yield (32). The candidate genes identified

in this study will be useful for understanding the underlying mechanism of this physiology.

Importantly then, we have identified many new genes responsible for high-yield, an economically most important trait. Most of these gene loci have not yet been functionally annotated, although a few of them belong to the β -expansin family or contain a zinc finger domain, which are known to play an important role in plant height, flower development, and light-regulated morphogenesis (33–35). We highlighted Os10g0555100, the knockout of which showed a different panicle architecture and a 23% reduction in height. We also note that our results suggest that the genes identified from cultivated lines in a pedigree could better reflect the real targets in plant breeding than the genes identified from artificial mutants. Our catalogue of 123 unannotated gene loci provides choices for downstream analysis. Our knockout and knockdown study of about half of these loci revealed that most of the genes in these loci are essential for rice phenotypes or for normal growth. Among the 159 genes we identified, there are at least 31 yield related genes, including 15 identified by knockout, 10 by knockdown and 6 previously reported. This proportion (19.5%) is significantly higher than the expectation (2.33 in 159 = 1.5%) based on the reported yield related genes in the rice genome ($p < 0.001$, $\chi^2 = 334$, $df = 1$, Chi-squared test with Yate's correction, see details in the *SI Appendix*, *SI Materials and Methods*). However, the alleles contributing to the Green Revolution are not necessarily null alleles, so our knockout and knockdown studies did not directly test the contribution of allelic changes to the Green Revolution. Gene replacements in IR8 or MH63 would directly reveal the contributions of the alleles, but IR8 and MH63 are difficult to transform and gene replacement is currently difficult in rice.

Our results also highlight clustering of unrelated genes with similar yield-associated phenotypes in the genome. This observation is of relevance to those hunting for complex trait genes and for those interested in genome evolution. For the former, it suggests that looking at neighbors of functionally relevant genes would be an effective way to look for functionally related genes. The clustering may reflect epistasis between genes or selection for co-expression. Previous QTL analysis also suggested that genes of similar phenotypic effects tend to cluster together (26), but this could also reflect allelic versions of control elements for a single gene. The fact that the knockouts of the clustered genes tend to have similar phenotypes suggests it is not the case.

Methods

Detailed materials and methods are outlined in *SI Appendix*, *SI Materials and Methods*

Plant materials and sequencing

The seeds of all rice accessions were obtained from International Rice Research Institute (IRRI) and China National Rice Research Institute (CN-RRRI) (Dataset S1). Pedigree information was obtained from the germplasm databases of IRRI and CNRRRI. All rice varieties were grown in the paddy-field. DNA samples were prepared from fresh leaves of a single plant using the Cetyltrimethyl Ammonium Bromide (CTAB) method and were sequenced at BGI-Shenzhen. Briefly, paired-end sequencing libraries with an insert size of ~500 bp were constructed for each plant, following the BGI-Shenzhen's instructions, and 2×100bp paired-end reads were generated on an Illumina HiSeq 2000. The sequencing reads of the 30 rice accessions have been deposited to the NCBI Short Read Archive under the Accession Numbers PRJNA271253 and SRR1060330. Indica cultivar 9311 callus RNA-seq data was downloaded from NCBI, BioProject PRJNA117345, SRR037711~SRR037724.

Construction of CRISPR Genome-Editing Vectors and target gene loci knockout

For each target locus, gRNAs were designed to target specific sites at the beginning of exons to cause a frame shift mutation. For each target, a pair of DNA oligonucleotides with appropriate cloning linkers were synthesized (BGI, Inc). Each pair of oligonucleotides were phosphorylated, annealed, and then ligated into BsaI-digested pRGEB31 vectors (addgene No.7722) (36). After transformation into *Escherichia coli* DH5- α , the resulting constructs were purified with Plasmid Mini kit (Genebase, Inc) for subsequent use in rice callus transformation. We selected Kasalath and Wuyungeng24 to be the background because they have a high transformation success rate while IR8 and MH63 are difficult to transform. Besides, Kasalath has a rather high

1225 stature and it is easy to observe when it becomes dwarf. Each construct was
1226 transformed into calli of Kasalath (an Indica) or Wuyugeng24 (a Japonica)
1227 by the method reported in a previous study (37). About ten transformed
1228 individuals were produced in two recipients for each vector (details in *SI*
1229 *Appendix*, Table S13, Table S14 and Table S15).

Genotype confirmation and phenotype observation

1230 The transgenic plants were examined under natural field conditions
1231 in the Experimental station of Nanjing University, Nanjing, China. For each
1232 plant, genomic DNA was extracted from fresh leaves by the CTAB method.
1233 In order to get double knockout mutants, we amplified the target region
1234 by PCR and confirmed the genotypes by Sanger sequencing. Primers were
1235 designed to make PCR products of ~1kb that contain the target sites. The
1236 results showed that 82.1% of transgenic plants had a knockout allele and
1237 79.5% had double knockout mutants. Phenotypes of the mutants were
1238 observed at different stages. Plant phenotypes were observed every three
1239 days to determine the changes in comparison with wild type rice plants. Plant
1240 height was measured after the stage of heading. Fertility and spike shape
1241 were observed when seeds were mature.

Footnotes

1242 ¹J.H.and J. L. contributed equally to this work.

1243 1. Mitchell-Olds T (2010) Complex-trait analysis in plants. *Genome Biol* 11(4):1–3.
1244 2. Sasaki T, Moore G eds. (1997) *Oryza: From Molecule to Plant* (Springer Netherlands, Dor-
1245 drecht) Available at: <http://link.springer.com/10.1007/978-94-011-5794-0> [Accessed October
1246 1, 2015].
1247 3. Zhu M (2007) Candidate Gene Identification Approach: Progress and Challenges. *Int J Biol*
1248 *Sci*:420–427.
1249 4. Korte A, Farlow A (2013) The advantages and limitations of trait analysis with GWAS: a
1250 review. *Plant Methods* 9(1):29.
1251 5. Mackay TFC (2001) The Genetic Architecture of Quantitative Traits. *Annu Rev Genet*
1252 35(1):303–339.
1253 6. Ott J, Kamatani Y, Lathrop M (2011) Family-based designs for genome-wide association
1254 studies. *Nat Rev Genet* 12(7):465–474.
1255 7. Pritchard JK, Stephens M, Donnelly P (2000) Inference of Population Structure Using
1256 Multilocus Genotype Data. *Genetics* 155(2):945–959.
1257 8. Brachi B, Morris GP, Borevitz JO (2011) Genome-wide association studies in plants: the
1258 missing heritability is in the field. *Genome Biol* 12(10):1–8.
1259 9. Kover PX, et al. (2009) A Multiparent Advanced Generation Inter-Cross to Fine-Map
1260 Quantitative Traits in Arabidopsis thaliana. *PLoS Genet* 5(7):c1000551.
1261 10. Shan Q, et al. (2013) Targeted genome modification of crop plants using a CRISPR-Cas
1262 system. *Nat Biotechnol* 31(8):686–688.
1263 11. Hazell PBR (2009) *The Asian Green Revolution* (Intl Food Policy Res Inst).
1264 12. Conway G (1998) *The Doubly Green Revolution: Food for All in the Twenty-first Century*
1265 (Comstock Pub. Associates).
1266 13. Peng S, Khush GS, Virk P, Tang Q, Zou Y (2008) Progress in ideotype breeding to increase
1267 rice yield potential. *Field Crops Res* 108(1):32–38.
1268 14. Khush GS (1999) Green revolution: preparing for the 21st century. *Genome* 42(4):646–655.
1269 15. Springer N (2010) Shaping a better rice plant. *Nat Genet* 42(6):475–476.
1270 16. Sasaki A, et al. (2002) Green revolution: A mutant gibberellin-synthesis gene in rice. *Nature*
1271 416(6882):701–702.
1272 17. Zhang J, et al. (2005) Features of the expressed sequences revealed by a large-scale analysis
1273 of ESTs from a normalized cDNA library of the elite indica rice cultivar Minghui 63. *Plant J*
1274 42(5):772–780.
1275 18. Wan X, et al. (2008) Quantitative Trait Loci (QTL) Analysis For Rice Grain Width and
1276 Fine Mapping of an Identified QTL Allele gw-5 in a Recombination Hotspot Region on
1277 Chromosome 5. *Genetics* 179(4):2239–2252.
1278 19. Li M, et al. (2011) Mutations in the F-box gene LARGER PANICLE improve the panicle
1279 architecture and enhance the grain yield in rice. *Plant Biotechnol J* 9(9):1002–1013.
1280 20. Xie W, et al. (2015) Breeding signatures of rice improvement revealed by a genomic variation
1281 map from a large germplasm collection. *Proc Natl Acad Sci* 112(39):E5411–E5419.
1282 21. Monna L, et al. (2002) Positional Cloning of Rice Semidwarfing Gene, sd-1: Rice “Green

1293 ²To whom correspondence should be addressed. E-mail:
1294 whli@uchicago.edu (W-H. L.) dtian@nju.edu.cn (D. T.) or bssldh@bath.ac.uk
1295 (L.D.H.)

1296 **Author contributions:** D.T, L.D.H. and W-H.L. designed the experiments
1297 and analyses. J.H. and J.L. organized all aspects of the project, J.H., J.L.,
1298 L.W. and S.Y. analyzed the sequence data. J.H., J.L. and J.Z. prepared plant
1299 materials and performed experimental confirmations. J.H., D.T., L.D.H. and
1300 W-H.L. wrote the paper.

The authors declare no conflict of interest.

Data deposition: Sequences have been deposited in the Sequence Read
1301 Archive, www.ncbi.nlm.nih.gov/sra (accession no. SRP051581).

Acknowledgements This work was supported by the National Major Special
1302 Project on New Varieties Cultivation for Transgenic Organisms (No.
1303 2016ZX08009001-003) and National Natural Science Foundation of China
1304 (91731308 and 31571267) to D.T. LDH is funded by ERC Advanced grant
1305 EvoGenMed (669207 — EvoGenMed — ERC-2014-ADG/ERC-2014-ADG). The
1306 authors thank the National Mid-term Genebank for Rice of China National
1307 Rice Research Institute and International Rice Research Institute for providing
1308 the rice germplasm collection.

1309 Revolution Gene” Encodes a Mutant Enzyme Involved in Gibberellin Synthesis. *DNA Res*
1310 9(1):11–17.
1311 22. Fang L, et al. (2012) Rolling-leaf14 is a 2OG-Fe (II) oxygenase family protein that modulates
1312 rice leaf rolling by affecting secondary cell wall formation in leaves. *Plant Biotechnol J*
1313 10(5):524–532.
1314 23. Zhou Y, et al. (2009) BC10, a DUF266-containing and Golgi-located type II membrane
1315 protein, is required for cell-wall biosynthesis in rice (*Oryza sativa* L.). *Plant J* 57(3):446–462.
1316 24. Nakashima K, et al. (2007) Functional analysis of a NAC-type transcription factor OsNAC6
1317 involved in abiotic and biotic stress-responsive gene expression in rice. *Plant J* 51(4):617–630.
1318 25. Hurst LD, Pál C, Lercher MJ (2004) The evolutionary dynamics of eukaryotic gene order.
1319 *Nat Rev Genet* 5(4):299–310.
1320 26. Cai H, Morishima H (2002) QTL clusters reflect character associations in wild and cultivated
1321 rice. *Theor Appl Genet* 104(8):1217–1228.
1322 27. Vazquez-Vilar M, et al. (2016) A modular toolbox for gRNA–Cas9 genome engineering in
1323 plants based on the GoldenBraid standard. *Plant Methods* 12. doi:10.1186/s13007-016-0101-2.
1324 28. Li JX, et al. (2000) Analyzing quantitative trait loci for yield using a vegetatively replicated
1325 F2 population from a cross between the parents of an elite rice hybrid. *Theor Appl Genet*
1326 101(1–2):248–254.
1327 29. Xing Y, et al. (2002) Characterization of the main effects, epistatic effects and their envi-
1328 ronmental interactions of QTLs on the genetic basis of yield traits in rice. *Theor Appl Genet*
1329 105(2–3):248–257.
1330 30. Yu SB, et al. (1997) Importance of epistasis as the genetic basis of heterosis in an elite rice
1331 hybrid. *Proc Natl Acad Sci* 94(17):9226–9231.
1332 31. Zhu Y, et al. (2012) Gene Discovery Using Mutagen-Induced Polymorphisms and Deep
1333 Sequencing: Application to Plant Disease Resistance. *Genetics* 192(1):139–146.
1334 32. Hedden P (2003) The genes of the Green Revolution. *Trends Genet* 19(1):5–9.
1335 33. Choi D, Lee Y, Cho H-T, Kende H (2003) Regulation of Expansin Gene Expression Affects
1336 Growth and Development in Transgenic Rice Plants. *Plant Cell Online* 15(6):1386–1398.
1337 34. Lee Y, Kende H (2001) Expression of β -Expansins Is Correlated with Internodal Elongation
1338 in Deepwater Rice. *Plant Physiol* 127(2):645–654.
1339 35. Takatsuiji H (1998) Zinc-finger transcription factors in plants. *Cell Mol Life Sci CMLS*
1340 54(6):582–596.
1341 36. Xie K, Yang Y (2013) RNA-Guided Genome Editing in Plants Using a CRISPR–Cas System.
1342 *Mol Plant* 6(6):1975–1983.
1343 37. Yang S, et al. (2013) Rapidly evolving R genes in diverse grass species confer resistance to
1344 rice blast disease. *Proc Natl Acad Sci* 110(46):18572–18577.
1345 38. Xie K, Yang Y (2013) RNA-Guided Genome Editing in Plants Using a CRISPR–Cas System.
1346 *Mol Plant* 6(6):1975–1983.