

# GpNet: Genomic Prediction Network Using Locally Connected Layers in Korean Native Cattle

**Hyo-Jun Lee**

Chungnam National University

**Dong Won Seo**

Chungnam National University

**Yoonji Chung**

Chungnam National University

**Doo Ho Lee**

Chungnam National University

**Yeung Kuk Kim**

Chungnam National University

**Jun Heon Lee**

Chungnam National University

**Hak Kyo Lee**

ChonBuk National University

**Cedric Gondro**

Michigan State University

**Young-Kuk Kim**

Chungnam National University

**Yeong Jun Koh**

Chungnam National University

**Seung Hwan Lee** (✉ [slee46@cnu.ac.kr](mailto:slee46@cnu.ac.kr))

Chungnam National University

---

## Research Article

**Keywords:** Genomic prediction, Deep learning, GWAS

**Posted Date:** June 29th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-622476/v1>

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

[Read Full License](#)

---



## RESEARCH

# GpNet: Genomic prediction network using locally connected layers in Korean native cattle

Hyo-Jun Lee<sup>1</sup>, Dong Won Seo<sup>2</sup>, Yoonji Chung<sup>2</sup>, Doo Ho Lee<sup>2</sup>, Yeung Kuk Kim<sup>2</sup>, Jun Heon Lee<sup>2</sup>, Hak Kyo Lee<sup>3</sup>, Cedric Gondro<sup>4</sup>, Young-Kuk Kim<sup>5</sup>, Yeong Jun Koh<sup>5\*</sup> and Seung Hwan Lee<sup>2\*</sup>

\*Correspondence: yjkoh@cnu.ac.kr;  
slee46@cnu.ac.kr

<sup>5</sup> Department of Computer  
Science & Engineering, Chungnam  
National University, 305-764,  
Daejeon, Korea

<sup>2</sup> Division of Animal and Dairy  
Science, Chungnam National  
University, 305-764, Daejeon,  
Korea

Full list of author information is  
available at the end of the article

## Abstract

**Background:** The use of DNA marker information for the prediction of genetic merit in animal and plant breeding, and susceptibility to disease in human medicine has become widespread. Therefore, an increasing number of methods have been proposed for more accurate and efficient genomic prediction. However, most of the commonly used models for genomic prediction only account for additive effects since most of them are designed based on the linear model.

**Results:** Here, we proposed a GpNet, a deep learning network for genomic prediction in Korean beef cattle. With a locally connected layer, GpNet can estimate LD-block effects of single nucleotide polymorphisms (SNP) with adjacent two or more SNPs closer to 3'-end. This operation is quite similar to how the DNA sequence is used in the translation process in which the RNA polymerase interprets DNA sequence by units of codons to downstream (3' to 5'). GpNet archived a superior performance than previous state-of-arts methods for beef carcass weight with a predictive ability of 0.721%. GpNet also found two significant quantitative trait locus (QTL) on the regions (bta 6:38464203–39816133, bta 14:25307116–29987025) for carcass weight. However, GpNet showed less performance than linear methods in backfat thickness and eye-muscle area.

**Conclusions:** GpNet outperformed the previous state-of-arts methods for beef carcass weight. However, GpNet cannot achieve superior performance in backfat thickness and eye-muscle area. We noticed that the lack of ability to estimate distant epistasis and dominance was the weakness of GpNet. Therefore, it remains a future research issue to expand GpNet to resolve these flaws and this further study will accelerate the new phase of the genomic prediction.

**Keywords:** Genomic prediction; Deep learning; GWAS

## 1 **Background**

2 The use of DNA marker information for the prediction of genetic merit in animal  
3 and plant breeding, and susceptibility to disease in human medicine has become  
4 widespread. This genomic information has been utilized primarily to detect regions  
5 of the genome that have an association with a specific phenotype (genome-wide  
6 association studies – GWAS) or to predict the genetic merit and phenotypes of  
7 individuals (genomic prediction) with many thousands of DNA markers, most com-  
8 monly single nucleotide polymorphisms (SNP), covering the entire genome. In hu-  
9 mans, genomic prediction has been widely used to predict disease risk and highly  
10 polygenic complex human traits [1, 2]. In agriculture, genomic prediction was used  
11 to estimate genomic breeding values (gEBV) which are then used to make selection  
12 decisions in a breeding population.

13 Most of the commonly used models for genomic prediction have been proposed  
14 based on the linear mixed models [3, 4]. Genomic best linear unbiased prediction  
15 (GBLUP) uses a mixed model approach which approximates a traditional infinites-  
16 imal model and assumes all SNP contribute a non-zero value to the genetic vari-  
17 ance [4]. It is a method that simply uses a genomic relationship matrix built from  
18 the genotypes instead of a traditional pedigree-based relationship matrix. Bayesian  
19 linear model assumes that some SNPs have zero effects, whereas others have small  
20 to moderate effects and uses the posterior distributions to the parameters of linear  
21 mixed model [3, 5]. Even though these methods showed the state-of-the-art perfor-  
22 mance on many populations, they only account for additive effects, since most of  
23 them are designed based on the linear model. Thus, extended methods to account  
24 for non-linearity effects, such as dominance and epistatic interactions, have been  
25 proposed recently [6, 7].

26 Deep learning is also a good alternative method to solve this problem. Recent ad-  
27 vances in deep neural networks have outperformed the state-of-the-art in computer  
28 vision, natural language processing, and audio recognition tasks [8, 9, 10, 11]. Using  
29 the local information of the input features, like image RGB-channel, text, or audio  
30 sequence, accelerated the successes of deep neural networks. Convolutional neural  
31 network (CNN), which is the most successful deep learning structure in computer  
32 vision, constitutes weights-shared filter operation for the adjacent region of input  
33 image [12]. Recurrent neural network (RNN) has been commonly used in sequence

34 to sequence problems, such as speech to text or natural language processing, gener-  
35 ating a new sequence of the specific time by using the information before that time  
36 of sequence [10]. These two networks hypothesize that the regions showing similar  
37 patterns in the input data could explain the similar features with each other. As  
38 shown in Fig 1(a), features in the image (*e.g.* hair, eye, nose, glass, and so on)  
39 have similar RGB-color patterns within the same features. Speech sound also has a  
40 similar frequency pattern with other similar sounds (Fig 1(b)).

**Figure 1** Example of image and sound data. (a) RGB-image; (b) Mel spectrogram of raw sound sequence.

41 Interestingly, the local information also can be addressed in the genomic predic-  
42 tion. The general concept of genomic prediction relies on the linkage disequilibrium  
43 (LD) between genetic markers and the unknown quantitative trait loci (QTL). With  
44 high-density SNP panels, the markers co-segregate with the causal mutations allow-  
45 ing their genetic effects to be indirectly estimated through the adjacent SNPs [3, 13].  
46 Considering this attribute of SNP data, genomic prediction model should estimate  
47 the effects of each LD-block consisting of locally adjacent two or more SNPs not a  
48 single SNP for the more accurate prediction. However, unlike the image and sound  
49 data, the LD-blocks even with the same SNP pattern do not always have the same  
50 effect on the individual traits. For the SNP data, it is more important to recognize  
51 how close each LD-block biologically to the unknown QTL than the SNP pattern.  
52 Therefore, a different approach from the previous deep learning networks, such as  
53 CNN or RNN, is required to use local information for genomic prediction. Prac-  
54 tically, the simple fully-connected networks that didn't use the local information  
55 usually showed better performance than other local-based networks in previous  
56 studies [14]. In addition, Zingaretti *et al.* [15] explored a convolutional neural net-  
57 work for genomic prediction of polyploid outcrossing species. Montesinos-López *et*  
58 *al.* [16] used deep learning to the multi-environment genomic prediction of plat  
59 complex traits. Pook *et al.* [17] applied the local convolutional neural networks on  
60 simulated maize and real Arabidopsis data. However, these studies can not provide  
61 clear evidences that deep neural networks can outperform the previous methods  
62 such as GBLUP or Bayesian linear models.

63 In this study, we proposed the Genomic prediction Network (GpNet) using a lo-  
 64 cally connected layer for genomic prediction in Korean native cattle. The locally  
 65 connected layer works similarly to the causal convolution, except that weights are  
 66 unshared, that is, a different set of weights is applied at each different LD-block. We  
 67 validated the performance of GpNet as follow processes. First, the GpNet perfor-  
 68 mances were evaluated on carcass weights, backfat thickness, and eye-muscle area of  
 69 Korean native cattle, and then its performance was compared with the GBLUP [4],  
 70 BayesA [3], and BayesLASSO [18]. Second, we also identified the candidate QTL  
 71 region using LD-block effects estimated by GpNet for each trait. Since there are  
 72 few results that deep learning outperformed the linear method, this study will be a  
 73 very interesting attempt in the field of genomic prediction.

## 74 Results

### 75 Model performance

76 Table 1 presents the performance of GpNet, GBLUP, BayesA and BayesLASSO. In  
 77 Table 1, we saw that GpNet (0.721) outperformed other linear method (GBLUP:  
 78 0.714, BayesA: 0.719, BayesLASSO: 0.712) in carcass weight (CWT). However,  
 79 GpNet showed less performance than the linear methods in backfat thickness (BF)  
 80 and eye-muscle area (EMA). Comparing between linear methods, BayesA showed  
 81 the best performance in BF (0.637) and EMA (0.728).

**Table 1 Predictive ability of each model.**

	CWT	BF	EMA
GpNet	<b>0.721 ± 0.018</b>	0.602 ± 0.01	0.708 ± 0.011
GBLUP	0.714 ± 0.018	0.626 ± 0.011	0.727 ± 0.014
BayesA	0.719 ± 0.019	<b>0.637 ± 0.012</b>	<b>0.728 ± 0.013</b>
BayesLASSO	0.712 ± 0.02	0.629 ± 0.011	0.723 ± 0.014

82 An LD-pruned SNP set was also used in this study. Briefly, pairs of SNPs in the  
 83 1000-kb with a squared correlation greater than 0.1 were noted, and these SNPs  
 84 were greedily pruned from the window until no such pairs remained. Finally, a total  
 85 of 21,629 SNPs was used as an LD-pruned SNP set. Table 2 shows the predictive  
 86 ability of each model with LD-pruned SNP. With the results using 50k SNP, GpNet  
 87 once again showed the best performance (0.712) in the CWT. However, in BF  
 88 and EMA, GpNet once again underperformed (BF: 0.589, EMA: 0.704) than linear  
 89 methods, and GBLUP showed the best performance (BF: 0.607, EMA: 0.72) for

90 these traits. Comparing with the 50K SNP results, the predictive abilities of all  
 91 models were decreased. These results corroborate the ideas of Manolio *et al.* [19],  
 92 who maintained that a marker subset may cause the missing heritability even though  
 the variants in subset can explain a large proportion of genetic variance.

**Table 2 Predictive ability of each model with LD-pruned SNP.**

	CWT	BF	EMA
GpNet	<b>0.712 ± 0.022</b>	0.589 ± 0.015	0.704 ± 0.013
GBLUP	0.701 ± 0.019	<b>0.607 ± 0.009</b>	<b>0.72 ± 0.016</b>
BayesA	0.709 ± 0.019	0.605 ± 0.01	0.718 ± 0.015
BayesLASSO	0.704 ± 0.019	0.606 ± 0.008	0.718 ± 0.015

93

#### 94 Identification of QTL using GpNet

95 To compare the QTL mapping of GpNet, we also estimated the SNPs effect using  
 96 the single marker linear mixed model (SMLMM). Fig 2 shows the mapping results  
 97 of two methods (GpNet and SMLMM). In the CWT (Fig 2(a)), two significant  
 98 peaks (bta 6:38464203–39816133 and bta 14:25307116-29987025) were found in both  
 99 methods and one peak (bta 4:4508164–4790444) was identified only by SMLMM.  
 100 All three regions were previously identified as QTLs for carcass weights in Korean  
 101 beef cattle [20, 21, 22, 23]. As shown in Fig 2(b), the genetic characteristic of BF  
 102 seemed to be more polygenic than CWT. Among the numerous significant peak,  
 103 we saw that the five loci (bta 2:99845066-107657675, bta 6:38464203–39816133, bta  
 104 13:53003864-54829615, bta 19:6241437-7833508, and bta 23:3320932-4781751) were  
 105 standing out at both methods. GpNet found a variant on bta 10:12412780 as the  
 106 most significant marker for BF. This variant is on the protein-coding region of  
 107 *DPP8* gene that was reported to attribute the dipeptidyl peptidase activity of cow  
 108 testis [24]. EMA seemed to have a similar genetic structure with CWT (Fig 2(c)).  
 109 These results seemed to be due to the genetic correlation between EMA and CWT.  
 110 In our data, EMA showed the 0.546 correlation with CWT. We can see that variant  
 111 on bta 27:23040097, which was not identified at CWT, was significant to EMA.  
 112 This variant is close to *DLC1*, which plays a key role in the regulation of small  
 113 GTP-binding proteins.

114 Table 3 shows the QTL region identified by both GpNet and SMLMM. In the  
 115 results, *SLIT2* seemed to be a key gene for complex traits of Korean beef cattle  
 116 since the association of *SLIT2* was replicated for CWT, BF, and EMA. *SLIT2* play

**Figure 2** Manhattan plots of SMLMM and GpNet for each trait. (a) carcass weight; (b) backfat thickness; (c) eye-muscle area; x-axis is SNP position and chromosome, y-axis is  $-\log_{10} P$ -value.

117 highly conserved roles in axon guidance and neuronal migration. A lot of genome-  
 118 wide studies have reported this gene to QTL for beef complex traits including organ  
 119 weight [25, 26], body weight [27], and fertility [28].

**Table 3** Predictive ability of each model with LD-pruned SNP.

Trait	Region	Close gene
CWT	4:4508164-4790444	<i>COBL</i>
	6:38464203-39816133	<b><i>SLIT2</i></b>
	14:25307116-29987025	<i>RAB2A</i>
BF	2:99845066-107657675	<i>SLC4A3</i>
	6:38464203-39816133	<b><i>SLIT2</i></b>
	13:53003864-54829615	<i>NTSR1</i>
	19:6241437-7833508	<i>C19H17orf67</i>
	23:3320932-4781751	<i>BMP5</i>
EMA	6:38464203-39816133	<b><i>SLIT2</i></b>
	10:20359799-21698931	<i>DHRS4</i>
	11:68687376-69190084	<i>LCLAT1</i>
	12:32687103-33576827	<i>LOC536660</i>
	14:25307116-29987025	<i>RAB2A</i>

## 120 Discussion

121 The flexibility of a deep neural network takes advantage of its non-linearity for ge-  
 122 nomic prediction in comparison to the traditional linear-based methods. Our pro-  
 123 posed GpNet can explore the epistasis of adjacent SNPs, called the LD-block effect  
 124 in this study, using a locally connected layer. This structure may stand to benefit  
 125 for the trait with some obvious causal loci, like CWT in this study. However, GpNet  
 126 under-performs than linear methods for BF and EMA, which have more polygenic  
 127 structure than CWT. In the polygenic trait, a lot of loci at various distances make  
 128 epistasis, which is quite challenging since GpNet only considers the interaction be-  
 129 tween adjacent SNPs. As shown in Fig 2, QTL mapping of GpNet for BF and EMA  
 130 was unclear, compared to SMLMM.

131 For capturing distant epistasis, a dilated convolution [29] would be a good alterna-  
 132 tive. It is equivalent to a convolution operation with a larger kernel filter by dilating  
 133 it with zeros. In the WaveNet [9], Stacked dilated convolution enabled networks to  
 134 have very large receptive fields for raw audio sequence. Otherwise, a dilated convo-



135 lution can estimate very large epistasis fields for genomic prediction. Even though a  
 136 locally connected layer included a dilated operation, it only consisted of a very low  
 137 delation rate due to the memory complexity  $O(n(d + 1))$ , where  $n$  is the number  
 138 of SNPs and  $d$  is depths of locally connected layer (Fig 3). On the other hand, di-  
 139 lated convolution effectively allows the networks to estimate epistasis with a much  
 140 large distance since it only requires the two shared parameters per layer (Fig 4).  
 141 Therefore, it remains a future work to incorporate dilated convolution to GpNet for  
 142 more accurate genomic prediction. Specifically, a locally connected layer at GpNet  
 143 can be used to estimate adjacent epistasis, while additional modules with dilated  
 144 convolution estimate distant epistasis. Then, the final epistasis can be accounted  
 145 by a combination of these two estimated values.

**Figure 3 Visualization of the locally connected layer.**  $n$  is the number of input SNPs,  $d$  means layer depths. Since weights are unshared at locally connected layer, the number of parameters at  $d$ -depth layer is  $(d + 1)n$

**Figure 4 Visualization of dilated convolution.**  $n$  is the number of input SNPs,  $d$  means layer depths. Since weights are shared at dilated convolution, the number of parameters at each layer is two.

146 Dominance also contributes to the total genetic potential for the phenotype  
 147 (Fig 5(a)). A nonlinear activation function, which is a critical part of the design  
 148 of a neural network, can allow such networks to compute nontrivial dominance.  
 149 GpNet adopts relu as a nonlinear activation function. However, relu is not suit-  
 150 able for identifying dominance since it is still linear for positive values. Instead of  
 151 relu, transformed sigmoid or tanh would be good options for calculating dominance  
 152 (Fig 5(b)).

**Figure 5 Type of dominance and candidate activations for dominance.** (a) Type of dominance; (b) Candidate activations for dominance.

153 Our proposed GpNet in this work did not consider the distant epistasis and dom-  
 154 inance. Fig 6 shows the identical model for genomic prediction. Even though some  
 155 candidate methods to implement each module (additive, epistasis, and dominance)  
 156 were discussed in this study, designing a full model in Fig 6 is of course still chal-

157 lenging. The optimal architectures (depth and width) of each module, which is  
158 dependent on the other modules' architectures, should be found separately and the  
159 aggregate methods to combine the features transformed from each module must  
160 be determined. In addition, the interaction of the three modules is also needed to  
161 be implemented in the full model. Therefore, it remains a future research issue to  
162 expand GpNet to an ideal genomic prediction model in Fig 6. This further research  
163 will accelerate the new phase for genomic prediction.

**Figure 6 Identical model and candidate methods for different type of genetic effects.**

## 164 Conclusions

165 In this paper, we presented a GpNet, a deep learning network for genomic predic-  
166 tion in Korean native cattle. With a locally connected layer, GpNet can estimate  
167 genetic effects of each LD-block consisting of neighboring two or more SNPs. In the  
168 results, the GpNet outperformed previous state-of-arts methods including GBLUP,  
169 BayesA, BayesLASSO in CWT. However, GpNet can't achieve superior perfor-  
170 mance than linear models in BF and EMA. Furthermore, GpNet did not consider  
171 the distant epistasis and dominance effects. To resolve these flaws, we discussed  
172 alternative methods in **Discussion** section. With these alternatives, it remains a  
173 future research issue to expand GpNet to an ideal genomic prediction model and  
174 this further research will accelerate the new phase for genomic prediction.

## 175 Materials and methods

### 176 Dataset

177 The commercial Korean native cattle population used in this study included 10000  
178 individuals (animals were born between 2010 ~ 2017 and samples were collected  
179 between 2013 ~ 2019) with phenotypic measurements for carcass weight (CWT/kg),  
180 eye-muscle area (EMA/cm<sup>2</sup>), and backfat thickness (BF/mm). BF and EMA were  
181 measured after a 24-hour chill at the junction between the 12th and 13th ribs.

182 Genomic DNA of the animals was extracted from longissimus thoracis muscle  
183 samples using a DNeasy Blood and Tissue Kit (Qiagen, Valencia, CA). A total of  
184 10000 samples were genotyped using the Illumina Bovine SNP50 BeadChip. SNP  
185 quality control was performed using PLINK1.9 software [18] based on the following

186 filtering criteria: minor allele frequency  $< 0.001$ ; SNP call rate  $< 0.1$ ; SNP on the  
 187 sex chromosomes. 1,853 SNP were excluded by the quality control filtering step  
 188 and the postfilter missing rate was 0.6% of the genotypes. These missing SNP were  
 189 then imputed with Eagle v2.4 [30] and a final total of 44,314 SNP were used in the  
 190 study. The dataset was split into train (80%), validation (10%), and test (10%) sets  
 191 to evaluate GpNet performance.

192 Notice that the National Institute of Animal Science (NIAS) in Rural Development  
 193 Administration (RDA) of South Korea approved the experimental procedures, and  
 194 all samples were taken under public animal health and welfare guidelines.

### 195 GpNet

196 In this paper, we proposed a new genomic prediction model operating on the  
 197 SNP data. We can write the SNP data set as a one-dimensional sequence  $\mathbf{x} =$   
 198  $\{x_1, x_2, \dots, x_n\}$  by base pair position. The goal of our proposed networks was to  
 199 assign an LD-block effect ( $LDB$ ) to SNP on the  $i$ -position ( $x_i$ ) as follows:

$$200 \quad LDB_{x_i} = \sum_{j=i-k}^i w_j^{(i)} x_j \quad (1)$$

201 where,  $x_j$  is genotype of SNP on the  $j$ -position;  $w_j^{(i)}$  is LD-effect of SNP  $x_j$  to  
 202 SNP  $x_i$ ;  $k$  is the LD-window size. Therefore, the assigned LD-block effect  $LDB_{x_i}$   
 203 is conditional estimated with a total of  $k$ -SNPs closer to 3'-end. This operation is  
 204 quite similar to how DNA sequence is used in the translation process in which the  
 205 RNA polymerase interprets DNA sequence by units of codons to downstream (3' to  
 206 5') and generates mRNA sequences with this information. To model this operation,  
 207 we opted a locally connected layer.

### 208 *Locally connected layer*

209 Fig 3 shows the visualization of the locally connected layer. A locally connected  
 210 layer was inspired by causal convolution [9] and local convolution [31]. By using a  
 211 locally connected layer, the network cannot violate the order of SNP. Otherwise,  
 212 The LD-block effect of SNP at the  $i$ -position cannot depend on any of the SNPs to  
 213 the 5'-end ( $x_{i+1}, x_{i+2}, \dots, x_n$ ). A locally connected layer adopts a dilated operation  
 214 where the window size is increased over layer depth. By this operation, a model can

215 estimate the LD effects from a wider LD-block as the depth of the networks get  
 216 deeper.

### 217 *Network Structure*

218 GpNet consists of the stacks of multiple locally connected layers (Fig 7). Both skip  
 219 connection [11] and relu activation [32] are used throughout the network to enable  
 220 training of a much deeper model. Given this structure, the model will estimate new  
 221 LD-block effects at each layer and then add them to the input SNPs. We scaled the  
 222 different layer depths  $d$  and stack number  $s$  for three different traits with a validation  
 223 set. In particular, we found the best values for CWT ( $d:4, s:1$ ), BF ( $d:4, s:3$ ) and  
 224 EMA ( $d:3, s:7$ ). Finally, from the last locally connected layer, one fully-connected  
 225 layer with non-activation (linear operation) yields a scalar, which is correspond to  
 226 the genomic estimated breeding value (gEBV).

**Figure 7 GpNet architecture.** LCL is locally connected layer;  $d$  and  $s$  is the layer depths and the number of stacks; *Relu* is relu activation; GpNet can be scaled with  $d$  and  $s$  for different traits.

### 227 Loss function

228 In the common deep learning approach, the model aims to predict a truth label  $y$   
 229 from input sequence  $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$ . However, the main purpose of genomic  
 230 prediction is to estimate a genetic portion of individual phenotype, not truth phe-  
 231 notype itself. Therefore, we added the  $h^2$  term in the mean squared error loss. Let  
 232  $\sigma_G^2$  and  $\sigma_P^2$  are genetic variance and phenotype variance. Then, the heritability can  
 233 be calculated by  $h^2 = \sigma_G^2 / \sigma_P^2$ . Finally, the loss function for training GpNet was  
 234 defined as follows:

$$235 \quad \mathcal{L}(y, \tilde{y}) = \frac{\sum_{p=1}^m (h^2 y_p - \tilde{y}_p)^2}{m} \quad (2)$$

236 where,  $y_p$  and  $\tilde{y}_p$  are observed phenotype and predicted phenotype of  $p$ -animal. By  
 237 shrinking  $y$  with  $h^2$ ,  $\tilde{y}$  will be converged to the gEBV at the end of model training.  
 238 The variance component,  $\sigma_G^2$  and  $\sigma_P^2$ , were estimated using an average information  
 239 restricted maximum likelihood [33] by implementing the AIREMLF90 program [34].  
 240 Table 4 shows the variance estimation results for each trait.

**Table 4** Variance component estimation results.

	CWT	BF	EMA
$\sigma_G^2$	962.3	9.5	52
$\sigma_P^2$	1495.6	15.5	90.2
$h^2$	0.392	0.378	0.366

#### 241 Implementation details

242 For training GpNet, we set a learning rate to  $10^{-4}$ . The training was iterated for  
 243 100 epochs with batch size 16 on GeForce RTX 3090. We employed Adam [35] op-  
 244 timizer to minimize the loss function. We determined network parameters, which  
 245 achieved the best performance on the validation set among all epochs. To validate  
 246 the GpNet ability, performances of GBLUP, BayesA and BayesLASSO were com-  
 247 pared with GpNet results. All test steps were repeated 5 times, with randomly split  
 248 train (80%), validation (10%), and test (10%) set. Since the linear methods don't  
 249 need the validation set, the 9000 animals (train set + validation set) were used  
 250 for training each linear model. We measured the model performance as the cor-  
 251 relation between phenotype and gEBV divided by the square root of heritability,  
 252  $cor(y, gEBV)/h$ , called predictive ability [36].

#### 253 Finding QTL with LD-block effects.

254 In addition to predict gEBV, GpNet also can be used for estimating the LD-block  
 255 effect of each SNP. As the layer feedforwards to the next layer, GpNet accumulates  
 256 the LD-effect of the SNP on the  $i$ -position ( $x_i$ ) as follows:

$$257 \quad f_d(x_i) = \sum_{j=i-k}^i w_j^{d(i)} f_{d-1}(x_j) + f_{d-1}(x_i) \quad (3)$$

258 where,  $f_d$  means the operation of the  $d$ -th layer;  $w_j^{d(i)}$  is LD-effect of SNP  $x_j$  to  
 259 SNP  $x_i$  on the  $d$ -th layer. Therefore, the final LD-block effects of SNP  $x_i$  can be  
 260 defined as:

$$261 \quad LD_{x_i} = f_E(x_i) - f_{E-1}(x_i) = \sum_{j=i-k}^i w_j^{E(i)} f_{E-1}(x_j) \quad (4)$$

262 where,  $f_E$  means the operation of the last locally connected layer;  $LD_{x_i}$  is final  
 263 LD-block effects of SNP  $x_i$ . Then, LD-block effect can be calculated by the differ-  
 264 ence between the outputs of the last layer  $f_E(x_i)$  and the one before it  $f_{E-1}(x_i)$ .

265 This LD-block effect, differing with the single SNP-effect estimated from the whole  
 266 population, can be different values for each individual since the LD-block effects are  
 267 estimated from individuals' SNPs patterns. To estimate the LD-block effects in this  
 268 study, we trained the whole population (10000 animals) to GpNet, and then 1000  
 269 animals with high-gEBV and low-gEBV were noted (Fig 8). We hypothesized that  
 270 the difference in gEBV ranking between these two groups (high and low) would be  
 271 reflected by the difference in LD-block effects of each individual. Therefore, we did  
 272 a t-test for finding a significant region as follows:

$$273 \quad Sig_i = \text{t-test}(H_{x_i}, L_{x_i}) \quad (5)$$

274 where  $Sig_i$  is significant value of  $i$ -position;  $H_{x_i} \in \mathbb{R}^{1000}$  and  $L_{x_i} \in \mathbb{R}^{1000}$  are LD-  
 275 block effects of SNP  $x_i$  in the high-gEBV group and low-gEBV group, separately.

**Figure 8 The process of QTL mapping.**

## 276 Abbreviations

- 277 • **SNP:** Single nucleotide polymorphisms
- 278 • **gEBV:** Genomic breeding values
- 279 • **GBLUP:** Genomic best linear unbiased prediction
- 280 • **CNN:** Convolutional neural network
- 281 • **RNN:** Recurrent neural network
- 282 • **LD:** Linkage disequilibrium
- 283 • **QTL:** Quantitative trait locus
- 284 • **CWT:** Carcass weight
- 285 • **BF:** Backfat thickness
- 286 • **EMA:** Eye-muscle area
- 287 • **SMLMM:** Single marker linear mixed model

## Declarations

Ethics approval and consent to participate

Notice that the National Institute of Animal Science (NIAS) in Rural Development Administration (RDA) of South Korea approved the experimental procedures, and all samples were taken under public animal health and welfare guidelines.

#### Availability of data and materials

All source code for this study is freely available for download at <https://github.com/gywns6287/GpNet>. Request for Genotype data can be made to Korea National Institute of Animal Science, Animal Genome & Bioinformatics Division (<http://www.nias.go.kr/english/sub/boardHtml.do?boardId=depintro>).

#### Competing interests

The authors declare that they have no competing interests.

#### Consent for publication

Not Applicable.

#### Funding

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (No. NRF-2019R1F1A1057605). CG was also supported by the National Institute of Food and Agriculture (AFRI Project No. 2019-67015-29323).

#### Authors' contributions

Conceptualization, HJL and YJK; Data curation, HKL and JHL; Methodology, DWS, YJC; Formal analysis, HJL and DHL; software, HJL, YKK, YKK; Writing – original draft, HJL; Writing – review & editing, SHL, CG. All authors have read and approved the final manuscript.

#### Acknowledgements

All data (phenotypes and genotypes) used in this study were provided by the Next Generation Biogreen21 project (PJ01316903), RDA. This study is part of an initiative to develop new genomic selection models using Artificial Intelligence and Machine Learning for Hanwoo cattle. We acknowledge the Next Generation Biogreen 21 Program, RDA.

#### Author details

<sup>1</sup> Department of Bio-AI Convergence, Chungnam National University, 305-764, Daejeon, Korea. <sup>2</sup> Division of Animal and Dairy Science, Chungnam National University, 305-764, Daejeon, Korea. <sup>3</sup> Department of Animal Biotechnology, ChonBuk National University, Jeonju, Korea. <sup>4</sup> Department of Animal Science, Michigan State University, East Lansing, MI, USA. <sup>5</sup> Department of Computer Science & Engineering, Chungnam National University, 305-764, Daejeon, Korea.

#### References

1. Abraham G, Inouye M. Genomic risk prediction of complex human disease and its clinical application [Journal Article]. *Current opinion in genetics & development*. 2015;33:10–16.
2. de Los Campos G, Vazquez AI, Fernando R, Klimentidis YC, Sorensen D. Prediction of complex human traits using the genomic best linear unbiased predictor [Journal Article]. *PLoS Genet*. 2013;9(7):e1003608.
3. Meuwissen TH, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide dense marker maps [Journal Article]. *Genetics*. 2001;157(4):1819–1829.
4. Hayes BJ, Visscher PM, Goddard ME. Increased accuracy of artificial selection by using the realized relationship matrix [Journal Article]. *Genetics research*. 2009;91(1):47–60.
5. Erbe M, Hayes B, Matukumalli L, Goswami S, Bowman P, Reich C, et al. Improving accuracy of genomic predictions within and between dairy cattle breeds with imputed high-density single nucleotide polymorphism panels [Journal Article]. *Journal of dairy science*. 2012;95(7):4114–4129.
6. Martini JW, Gao N, Cardoso DF, Wimmer V, Erbe M, Cantet RJ, et al. Genomic prediction with epistasis models: on the marker-coding-dependent performance of the extended GBLUP and properties of the categorical epistasis model (CE) [Journal Article]. *BMC bioinformatics*. 2017;18(1):1–16.
7. Da Y, Wang C, Wang S, Hu G. Mixed model methods for genomic prediction and variance component estimation of additive and dominance effects using SNP markers [Journal Article]. *PLoS one*. 2014;9(1):e87666.

8. Tan M, Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In: International Conference on Machine Learning. PMLR; 2019. p. 6105–6114.
9. Oord Avd, Dieleman S, Zen H, Simonyan K, Vinyals O, Graves A, et al. Wavenet: A generative model for raw audio [Journal Article]. arXiv preprint arXiv:160903499. 2016;.
10. Cho K, Van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation [Journal Article]. arXiv preprint arXiv:14061078. 2014;.
11. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition; 2016. p. 770–778.
12. Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks [Journal Article]. Advances in neural information processing systems. 2012;25:1097–1105.
13. Kizilkaya K, Fernando R, Garrick D. Genomic prediction of simulated multibreed and purebred performance using observed fifty thousand single nucleotide polymorphism genotypes [Journal Article]. Journal of animal science. 2010;88(2):544–551.
14. Bellot P, de Los Campos G, Pérez-Enciso M. Can deep learning improve genomic prediction of complex human traits? [Journal Article]. Genetics. 2018;210(3):809–819.
15. Zingaretti LM, Gezan SA, Ferrão LFV, Osorio LF, Monfort A, Muñoz PR, et al. Exploring deep learning for complex trait genomic prediction in polyploid outcrossing species. Frontiers in plant science. 2020;11:25.
16. Montesinos-López A, Montesinos-López OA, Gianola D, Crossa J, Hernández-Suárez CM. Multi-environment genomic prediction of plant traits using deep learners with dense architecture. G3: Genes, Genomes, Genetics. 2018;8(12):3813–3828.
17. Pook T, Freudenthal J, Korte A, Simianer H. Using local convolutional neural networks for genomic prediction [Journal Article]. Frontiers in genetics. 2020;11.
18. Park T, Casella G. The bayesian lasso [Journal Article]. Journal of the American Statistical Association. 2008;103(482):681–686.
19. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. Nature. 2009;461(7265):747–753.
20. Li Y, Lee JH, Lee YM, Kim JJ. Application of linkage disequilibrium mapping methods to detect QTL for carcass quality on chromosome 6 using a high density SNP map in Hanwoo [Journal Article]. Asian-Australasian Journal of Animal Sciences. 2011;24(4):457–462.
21. Lee SH, Choi BH, Lim D, Gondro C, Cho YM, Dang CG, et al. Genome-wide association study identifies major loci for carcass weight on BTA14 in Hanwoo (Korean cattle) [Journal Article]. PLoS One. 2013;8(10):e74677.
22. Sudrajat P, Sharma A, Dang CG, Kim JJ, Kim KS, Lee JH, et al. Validation of single nucleotide polymorphisms associated with carcass traits in a commercial Hanwoo population. Asian-Australasian journal of animal sciences. 2016;29(11):1541.
23. Lee HJ, Chung YJ, Jang S, Seo DW, Lee HK, Yoon D, et al. Genome-wide identification of major genes and genomic prediction using high-density and text-mined gene-based SNP panels in Hanwoo (Korean cattle). PloS one. 2020;15(12):e0241848.
24. Dubois V, Ginneken CV, Cock HD, Lambeir AM, Veken PVd, Augustyns K, et al. Enzyme activity and immunohistochemical localization of dipeptidyl peptidase 8 and 9 in male reproductive tissues [Journal Article]. Journal of Histochemistry & Cytochemistry. 2009;57(6):531–541.
25. An B, Xia J, Chang T, Wang X, Miao J, Xu L, et al. Genome-wide association study identifies loci and candidate genes for internal organ weights in Simmental beef cattle. Physiological genomics. 2018;50(7):523–531.
26. Raza SHA, Khan S, Amjadi M, Abdelnour SA, Ohran H, Alanazi KM, et al. Genome-wide association studies reveal novel loci associated with carcass and body measures in beef cattle. Archives of Biochemistry and Biophysics. 2020;p. 108543.
27. Smith JL, Wilson ML, Nilson SM, Rowan TN, Oldeschulte DL, Schnabel RD, et al. Genome-wide association and genotype by environment interactions for growth traits in US Gelbvieh cattle. BMC genomics. 2019;20(1):1–13.
28. Höglund JK, Buitenhuis B, Gulbrandsen B, Lund MS, Sahana G. Genome-wide association study for female



- fertility in Nordic Red cattle. *BMC genetics*. 2015;16(1):1–11.
29. Chen LC, Papandreou G, Kokkinos I, Murphy K, Yuille AL. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*. 2017;40(4):834–848.
  30. Loh PR, Danecek P, Palamara PF, Fuchsberger C, Reshef YA, Finucane HK, et al. Reference-based phasing using the Haplotype Reference Consortium panel [Journal Article]. *Nature genetics*. 2016;48(11):1443.
  31. Taigman Y, Yang M, Ranzato M, Wolf L. Deepface: Closing the gap to human-level performance in face verification. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*; 2014. p. 1701–1708.
  32. Nair V, Hinton GE. Rectified linear units improve restricted boltzmann machines. In: *Icml*; 2010. .
  33. Meyer K. An 'average information' restricted maximum likelihood algorithm for estimating reduced rank genetic covariance matrices or covariance functions for animal models with equal design matrices [Journal Article]. *Genetics Selection Evolution*. 1997;29(2):1–20.
  34. Misztal I, Tsuruta S, Strabel T, Auvray B, Druet T, Lee D, et al. BLUPF90 and related programs (BGF90). In: *Proceedings of the 7th world congress on genetics applied to livestock production*. vol. 33; 2002. p. 743–744.
  35. Kingma DP, Ba JL. Adam: A method for stochastic gradient descent. In: *ICLR: International Conference on Learning Representations*; 2015. p. 1–15.
  36. Legarra A, Robert-Granié C, Manfredi E, Elsen JM. Performance of genomic selection in mice [Journal Article]. *Genetics*. 2008;180(1):611–618.

# Figures

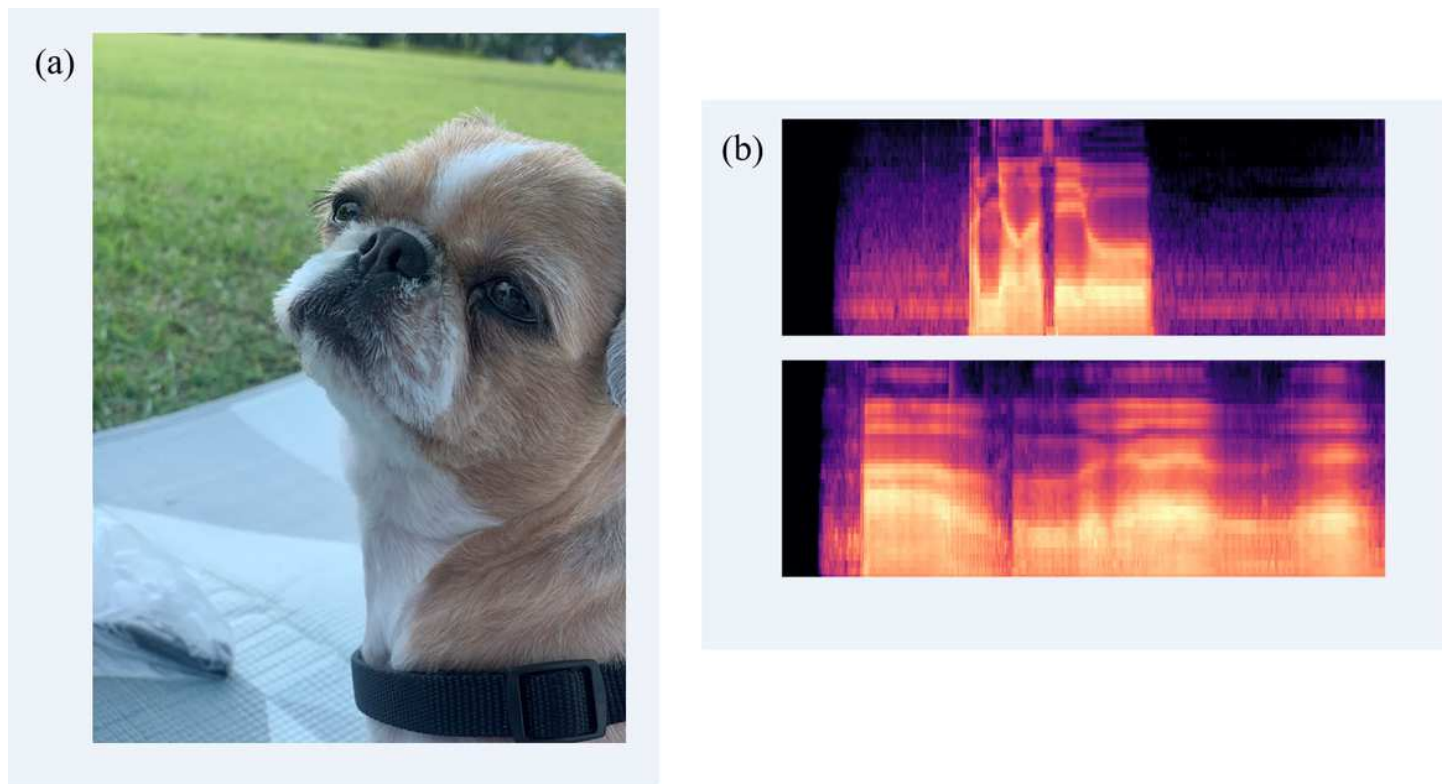


Figure 1

Example of image and sound data. (a) RGB-image; (b) Mel spectrogram of raw sound sequence.

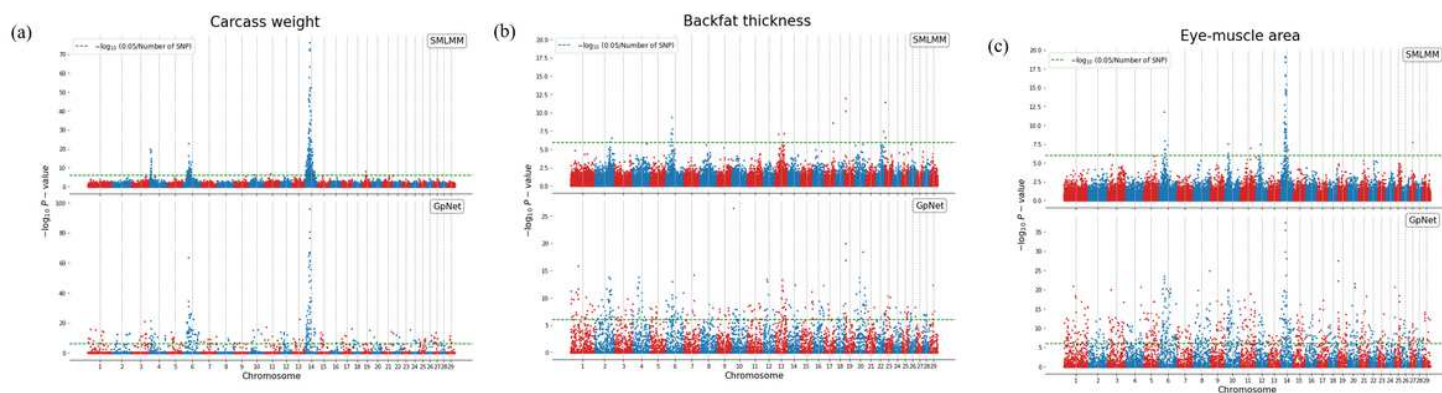
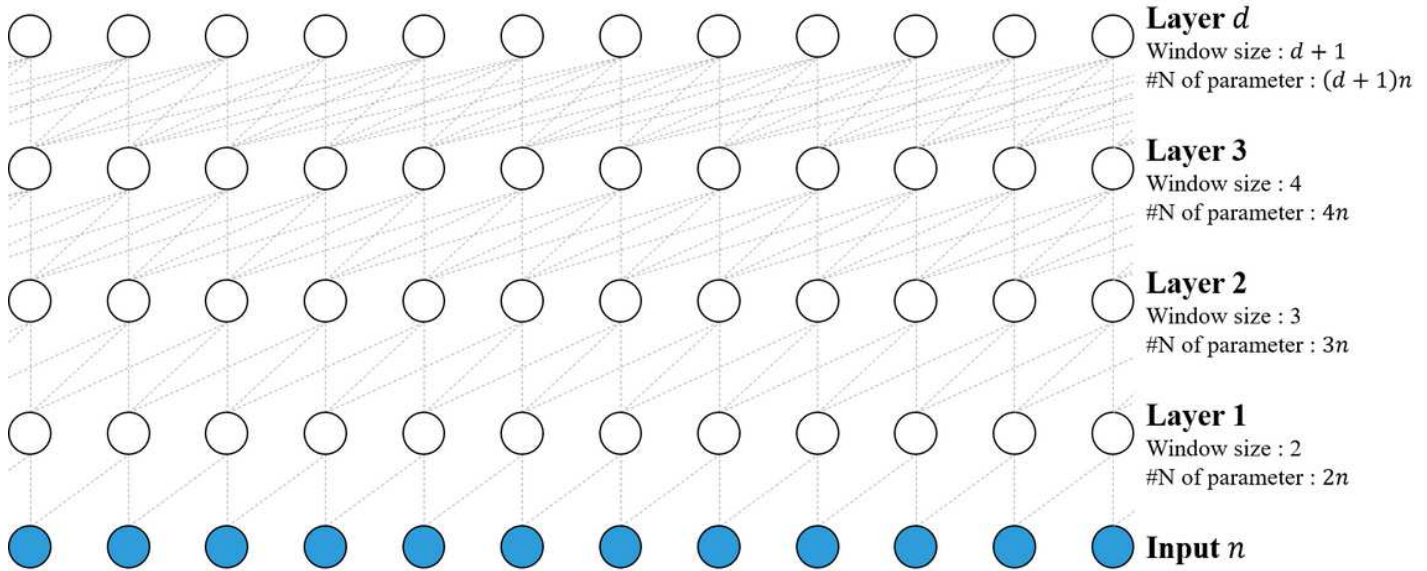


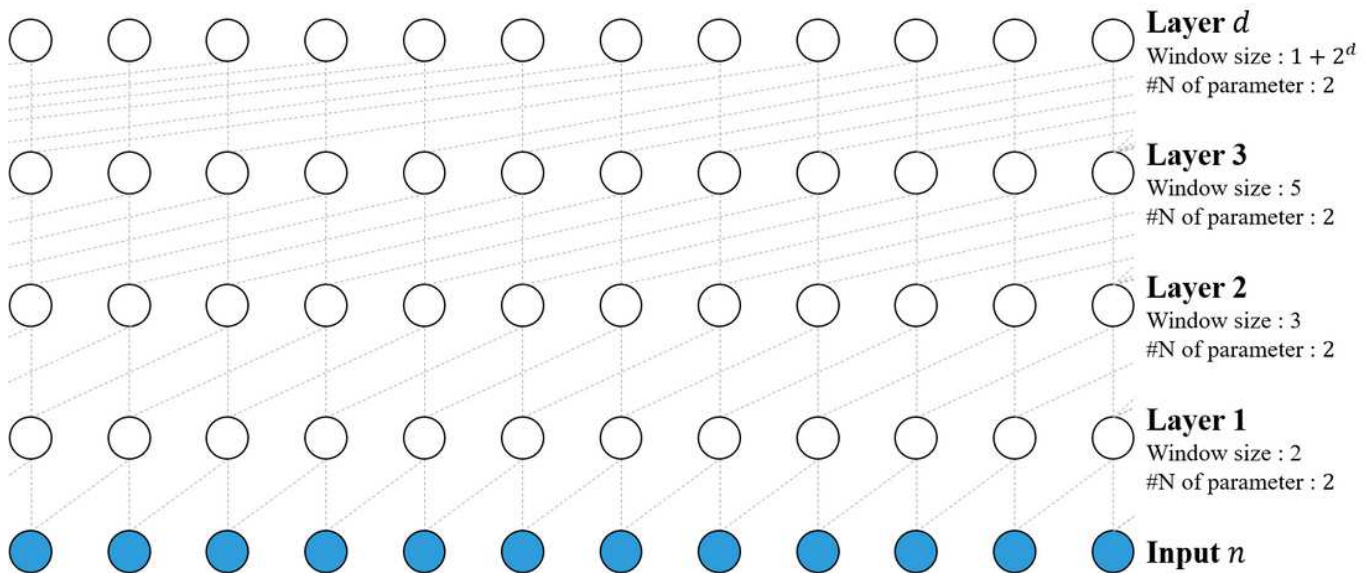
Figure 2

Manhattan plots of SMLMM and GpNet for each trait. (a) carcass weight; (b) backfat thickness; (c) eye-muscle area; x-axis is SNP position and chromosome, y-axis is  $-\log_{10} P$ -value



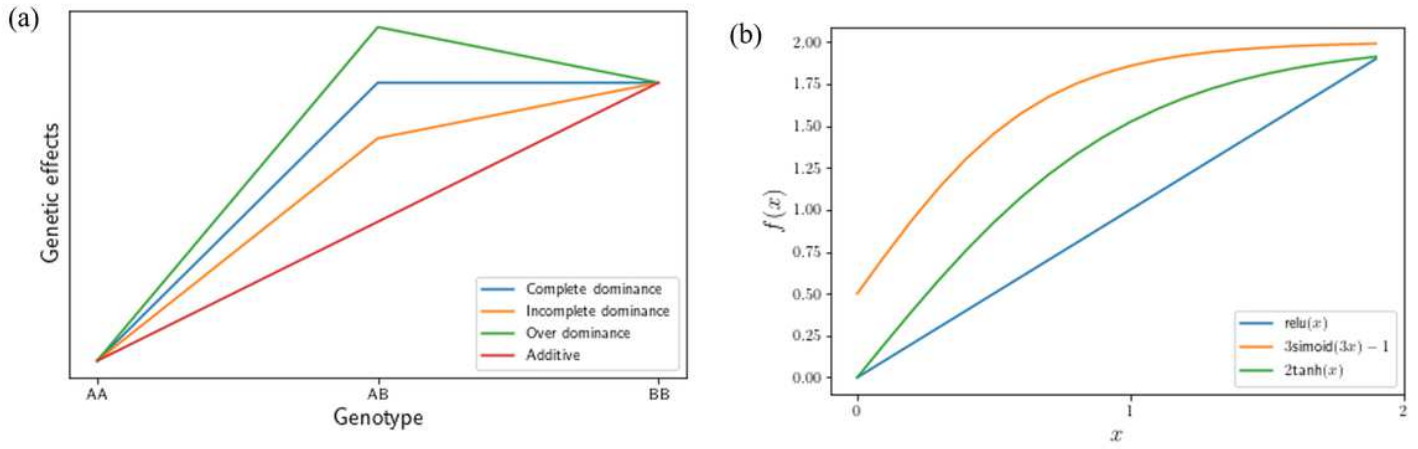
**Figure 3**

Visualization of the locally connected layer.  $n$  is the number of input SNPs,  $d$  means layer depths. Since weights are unshared at locally connected layer, the number of parameters at  $d$ -depth layer is  $(d + 1)n$



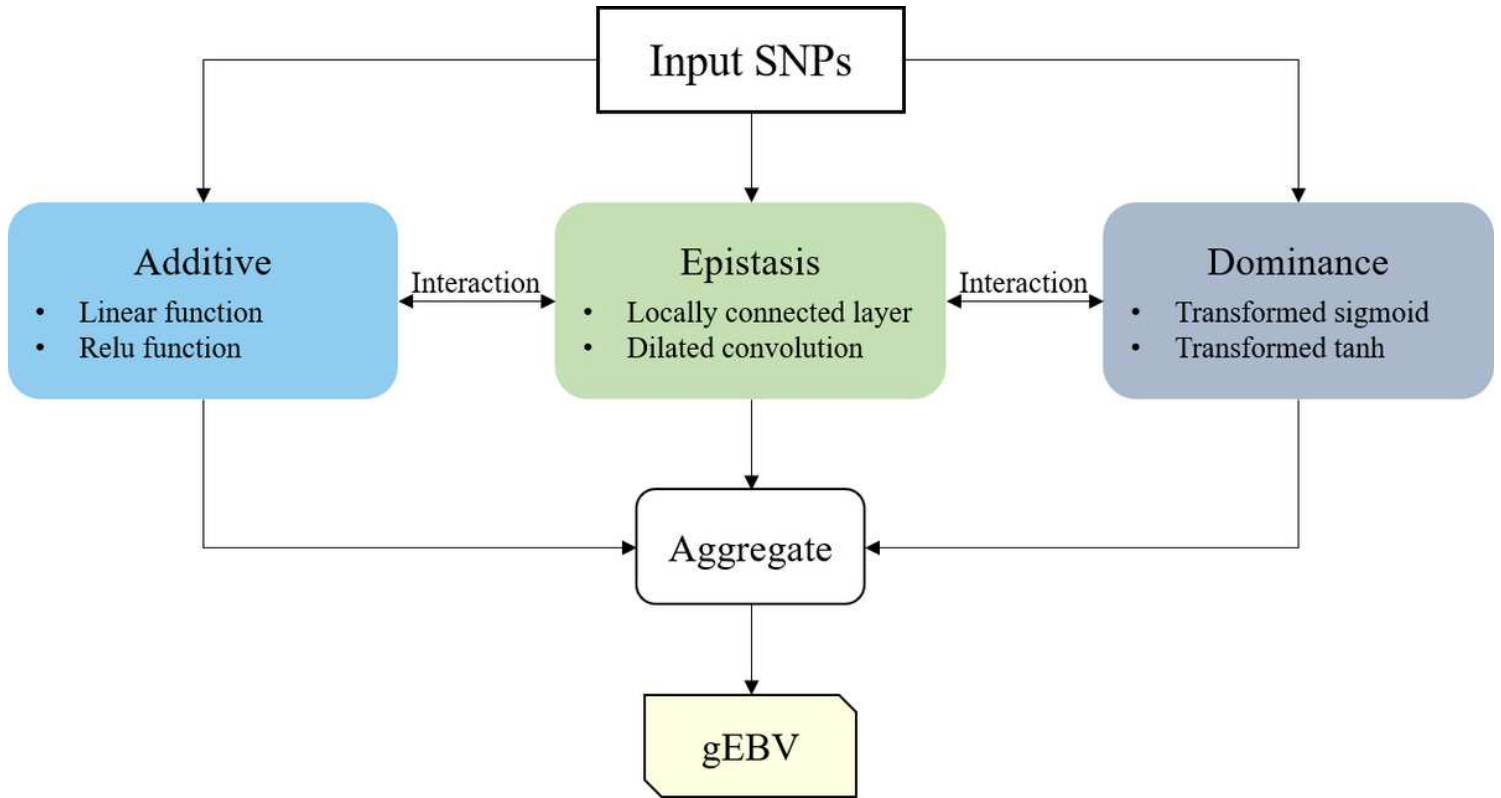
**Figure 4**

Visualization of dilated convolution.  $n$  is the number of input SNPs,  $d$  means layer depths. Since weights are shared at dilated convolution, the number of parameters at each layer is two.



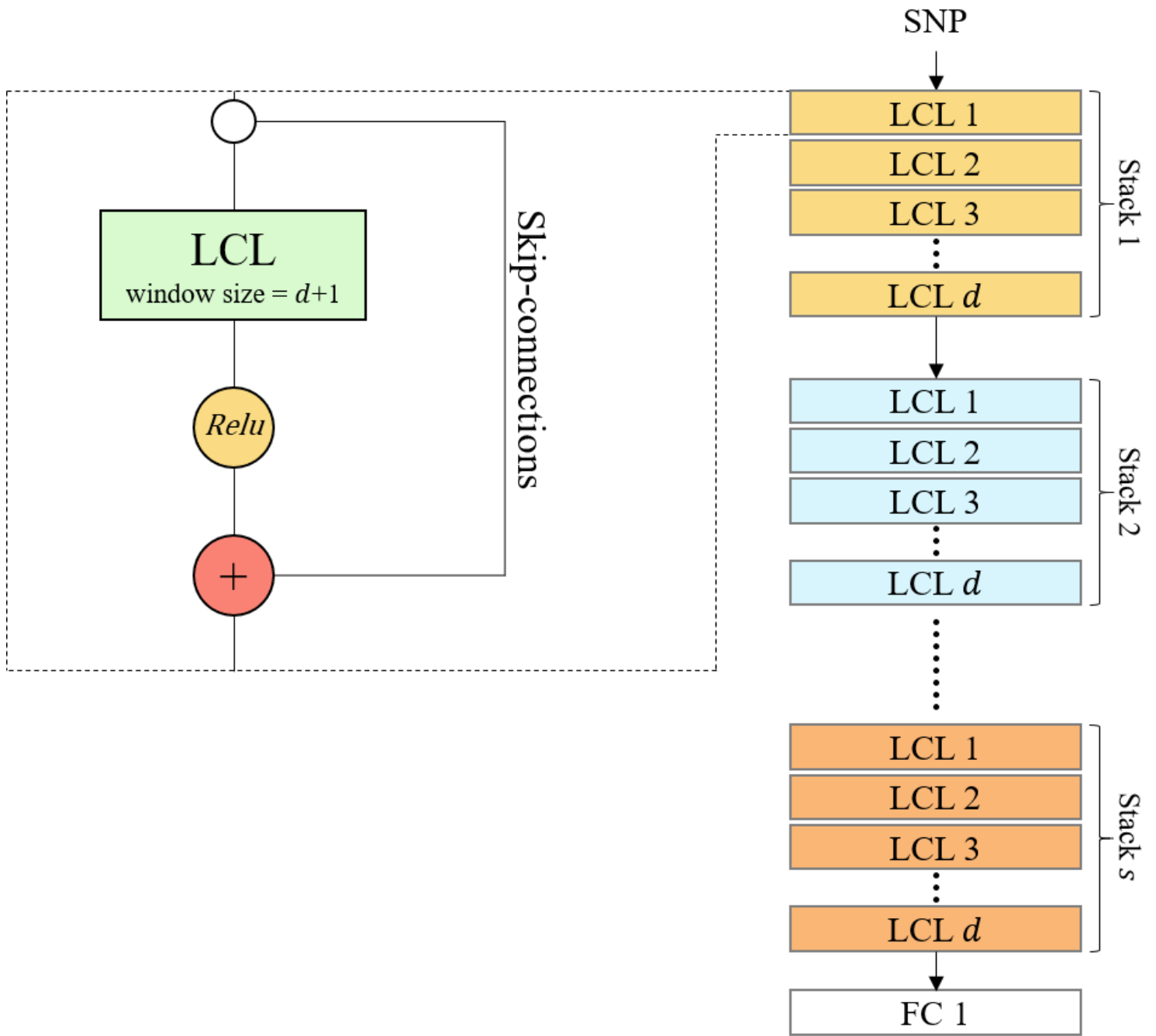
**Figure 5**

Type of dominance and candidate activations for dominance. (a) Type of dominance; (b) Candidate activations for dominance.



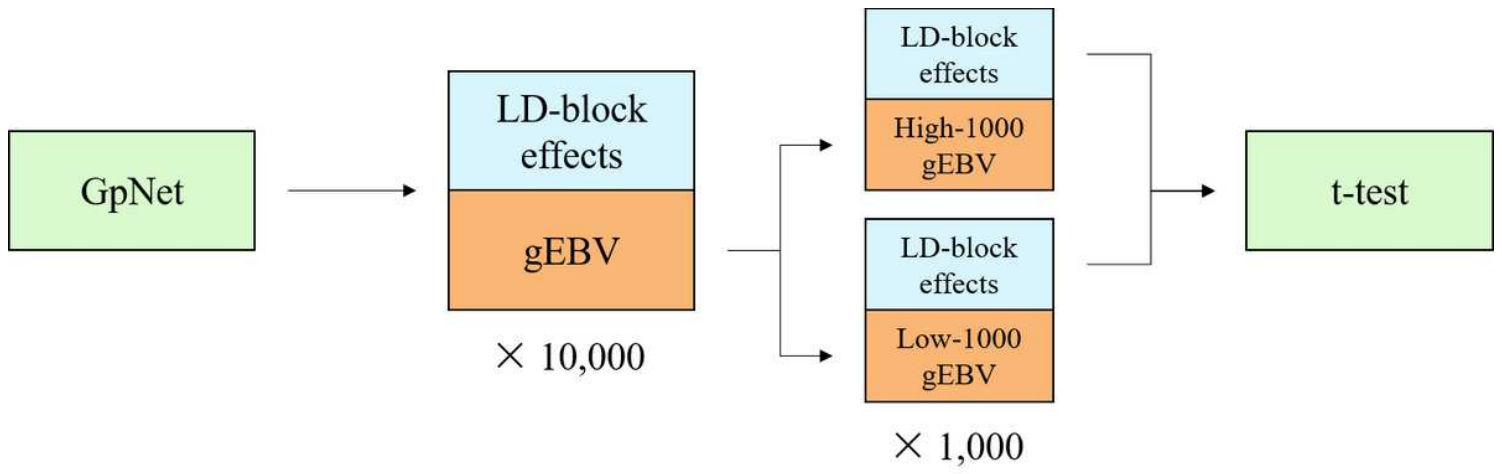
**Figure 6**

Identical model and candidate methods for different type of genetic effects.



**Figure 7**

GpNet architecture. LCL is locally connected layer;  $d$  and  $s$  is the layer depths and the number of stacks; Relu is relu activation; GpNet can be scaled with  $d$  and  $s$  for different traits



**Figure 8**

The process of QTL mapping.