

Cluster Analysis of a Symbolic Regression Search Space



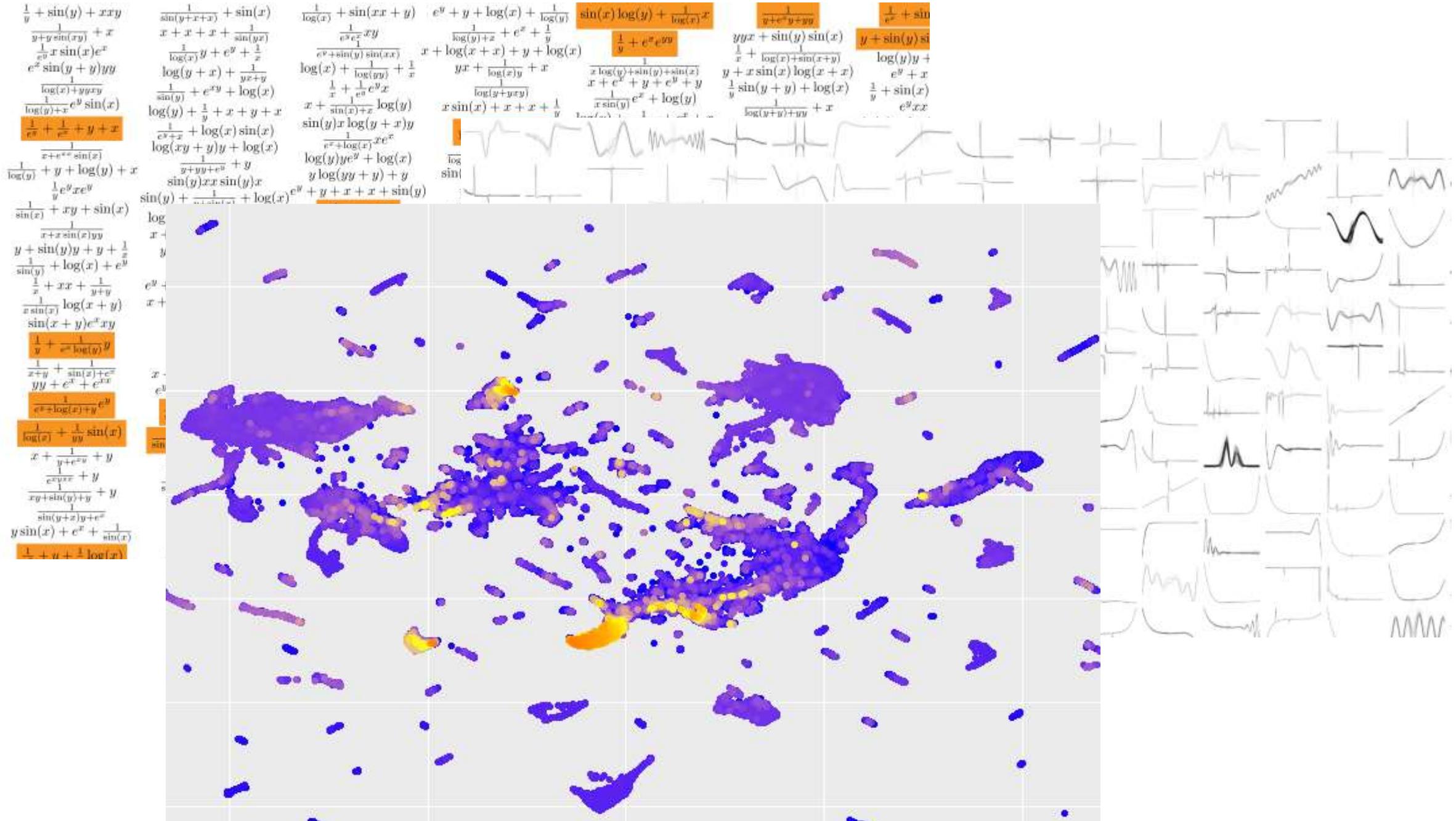
G. Kronberger, L. Kammerer, B. Burlacu, S. Winkler, M. Kommenda, M. Affenzeller



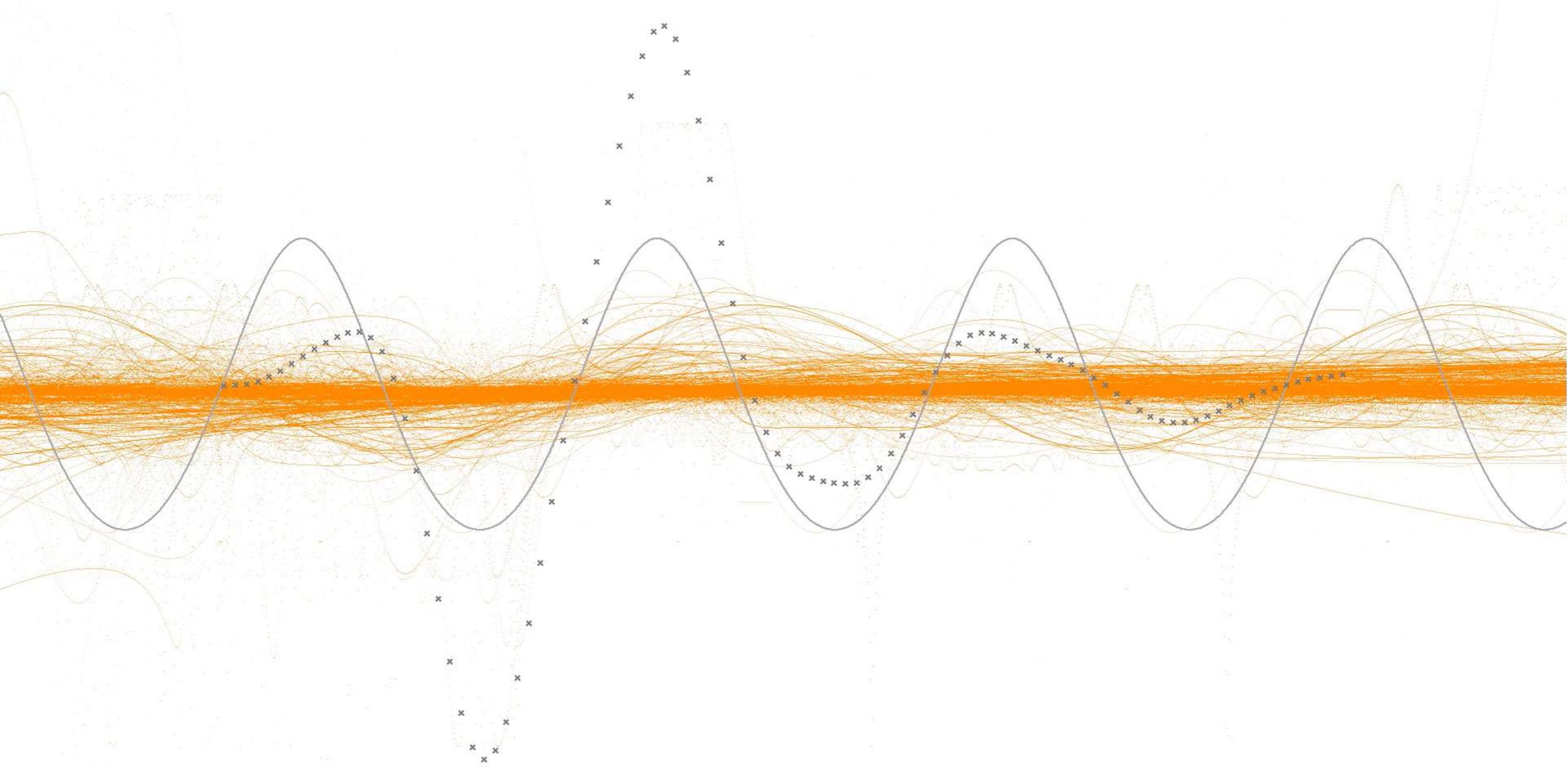
SymReg

JOSEF RESSEL CENTER FOR
SYMBOLIC REGRESSION

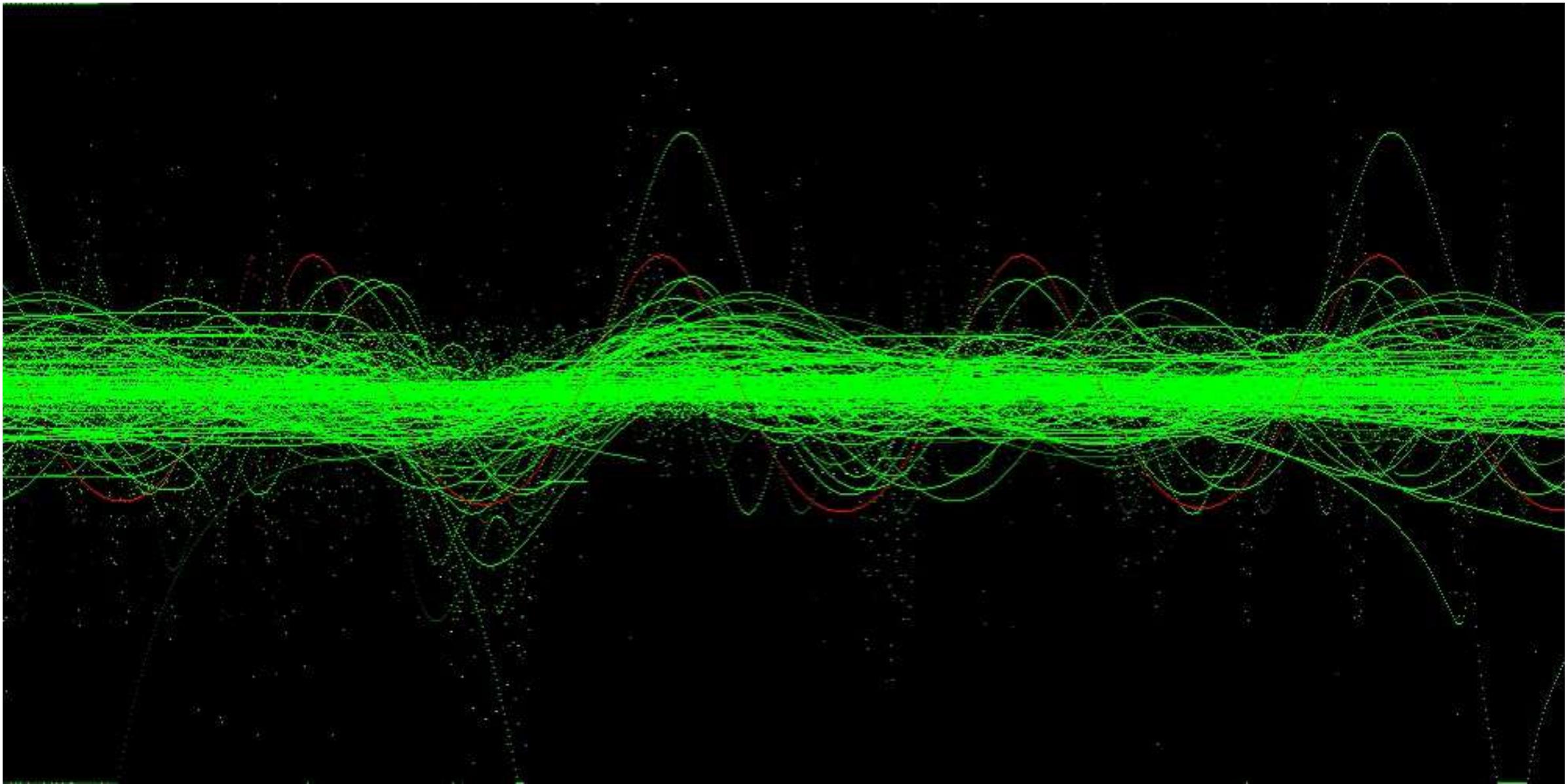




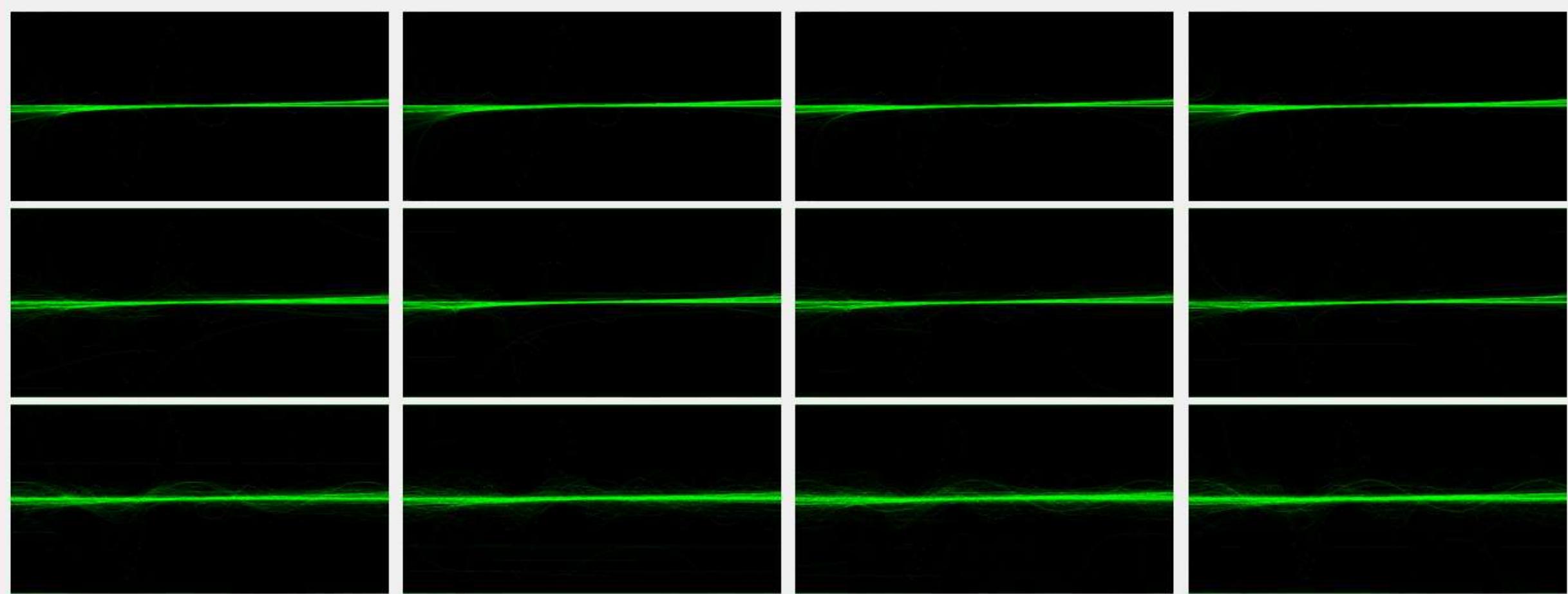
<https://youtu.be/IGqfAx7igkQ>



`EXP(COS(SIN((-1*X) + COS(COS(SIN(SIN(COS(LOG(((NaN*X) + (NaN)) / ((-1*X) + 6.3)))))))))))`



<https://youtu.be/yKh8ux68kMc>



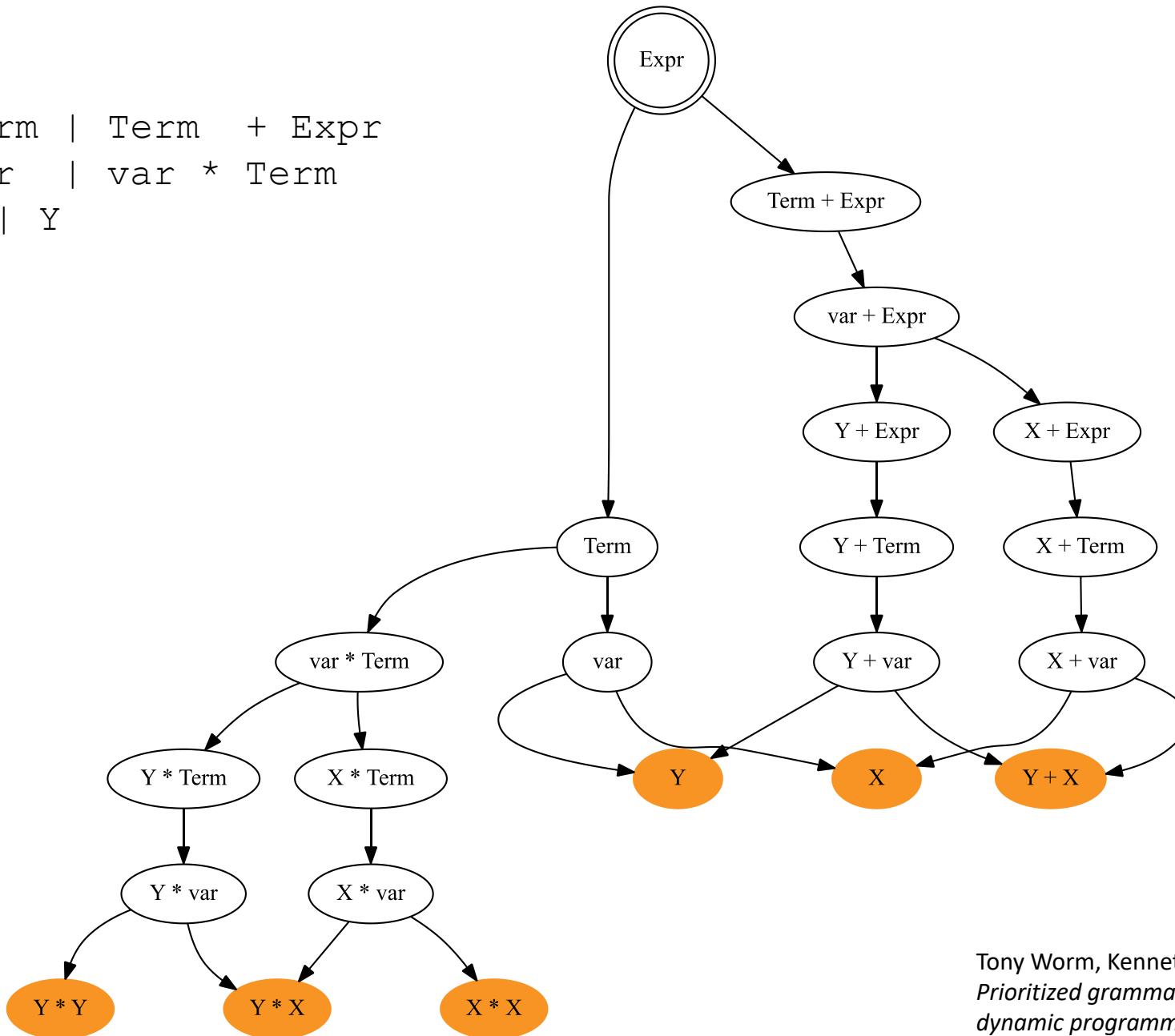
<https://youtu.be/9OStdHbTsqk>

$G(\text{Expr}) :$

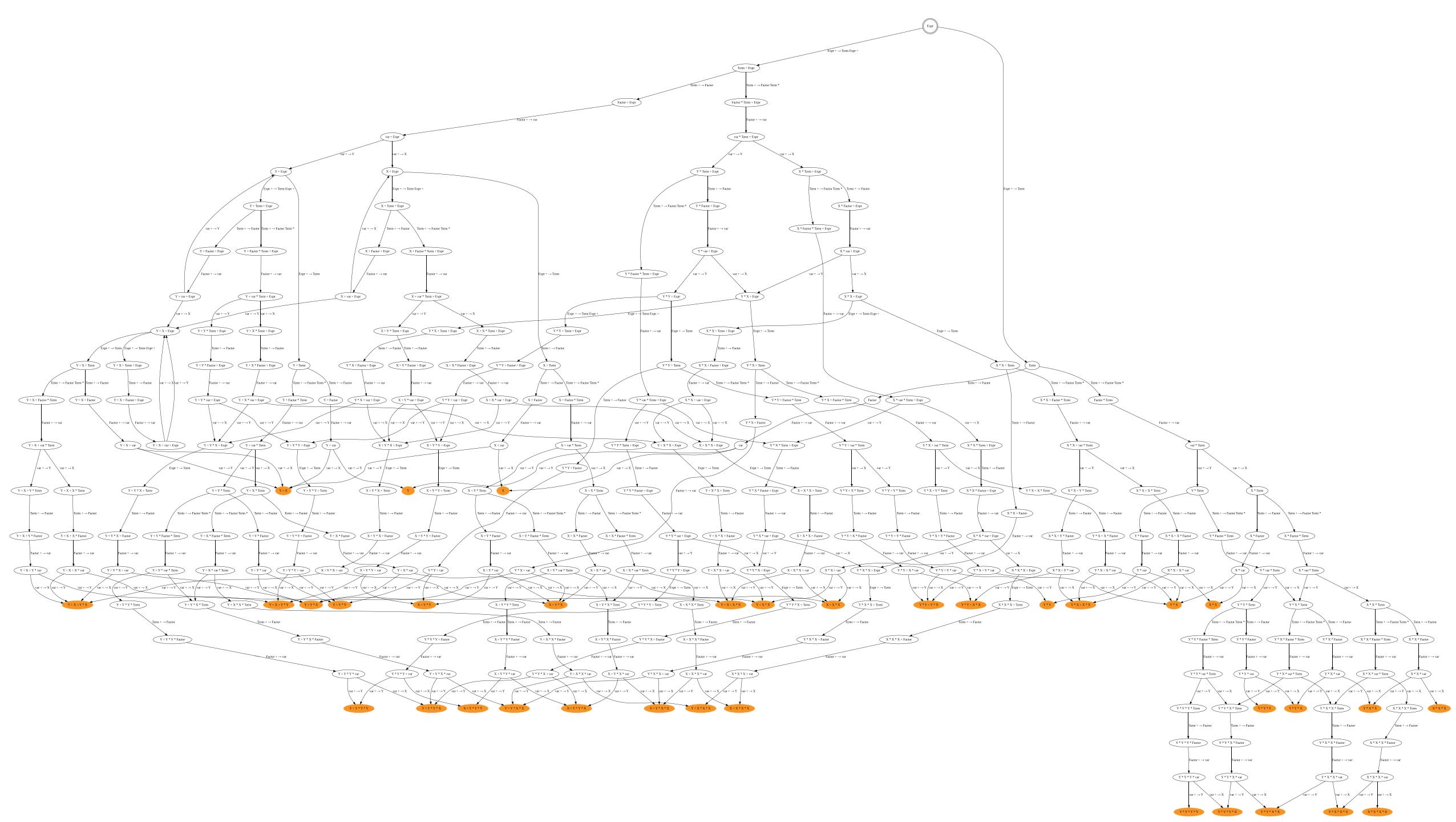
$\text{Expr} \rightarrow \text{Term} \mid \text{Term} + \text{Expr}$

$\text{Term} \rightarrow \text{var} \mid \text{var} * \text{Term}$

$\text{var} \rightarrow \text{X} \mid \text{Y}$



Tony Worm, Kenneth Chiu:
Prioritized grammar enumeration: symbolic regression by dynamic programming. GECCO 2013: 1021-1028



G(Expr) :

Expr → Term "+" Expr | Term

Term → Factor "*" Term | Factor | "1/(" InvExpr ")"

Factor → VarFac | ExpFac | LogFac | SinFac

VarFac → <variable>

ExpFac → "exp(" SimpleTerm ")"

LogFac → "log(" SimpleExpr ")"

SinFac → "sin(" SimpleExpr ")"

SimpleExpr → SimpleTerm "+" SimpleExpr | SimpleTerm

SimpleTerm → VarFac "*" SimpleTerm | VarFac

InvExpr → InvTerm "*" InvExpr | InvTerm

InvTerm → Factor "*" InvTerm | Factor

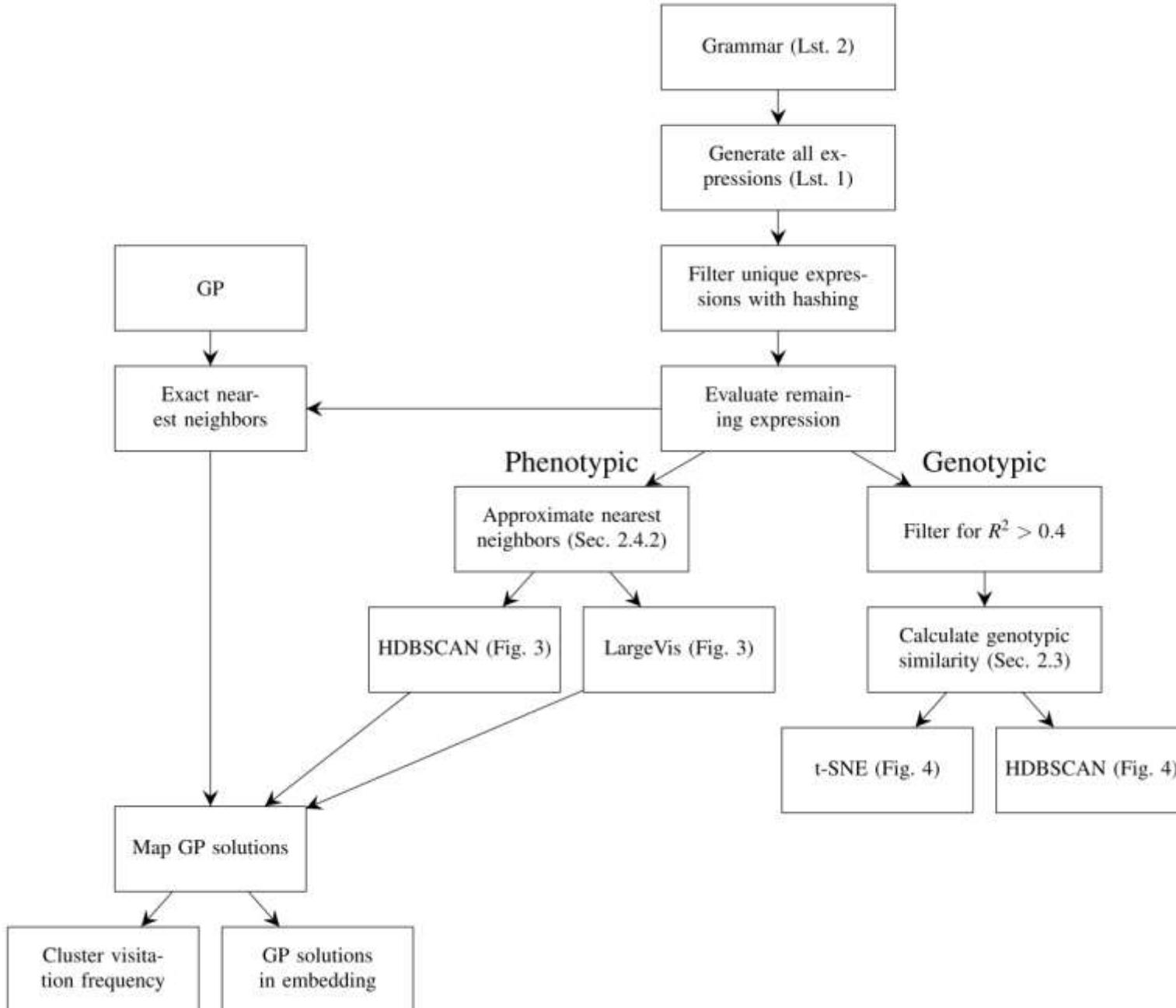
$$\frac{1}{x}+\sin(y)+x^2y$$

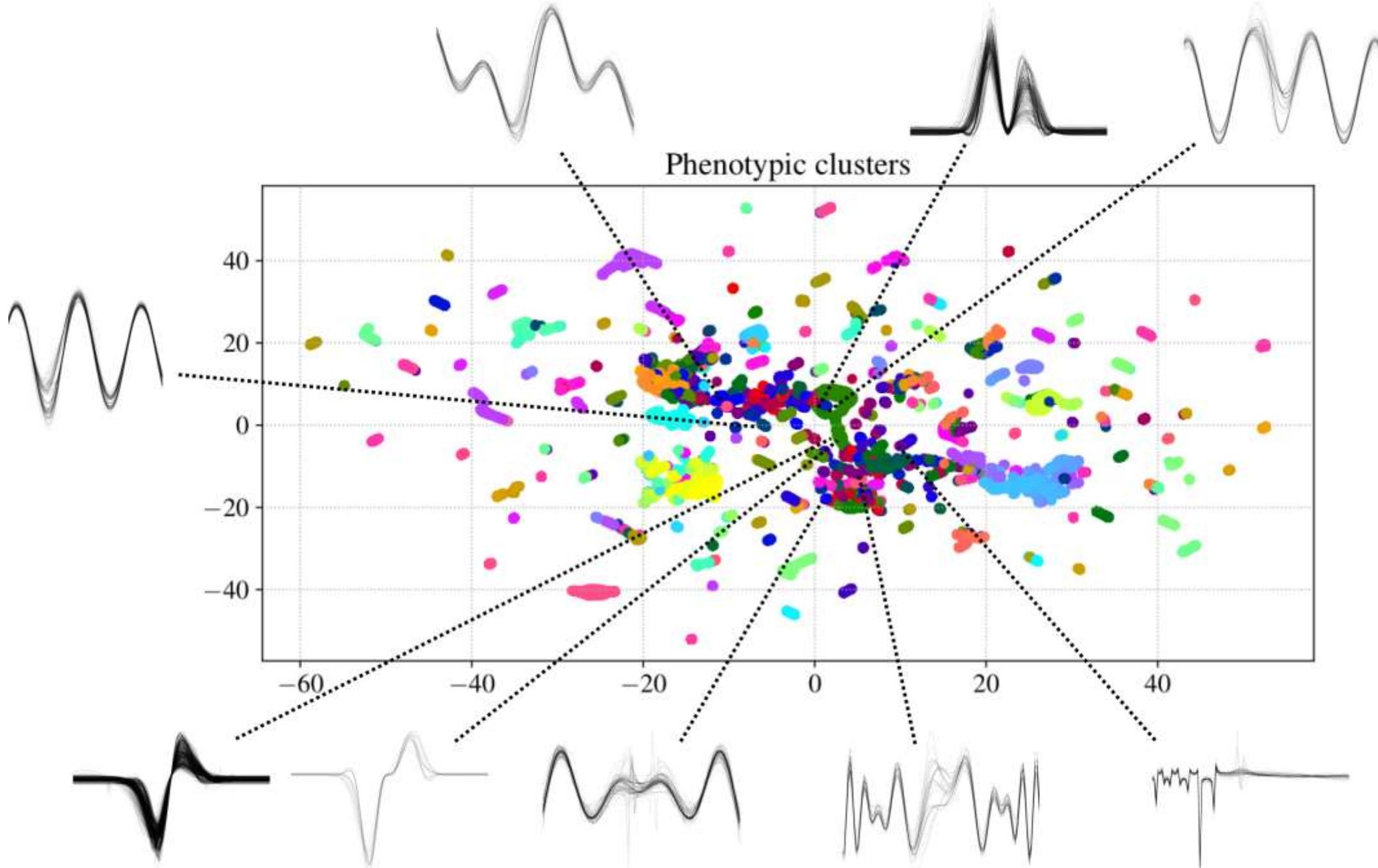
$$\frac{1}{e^y}x\sin(x)e^x$$

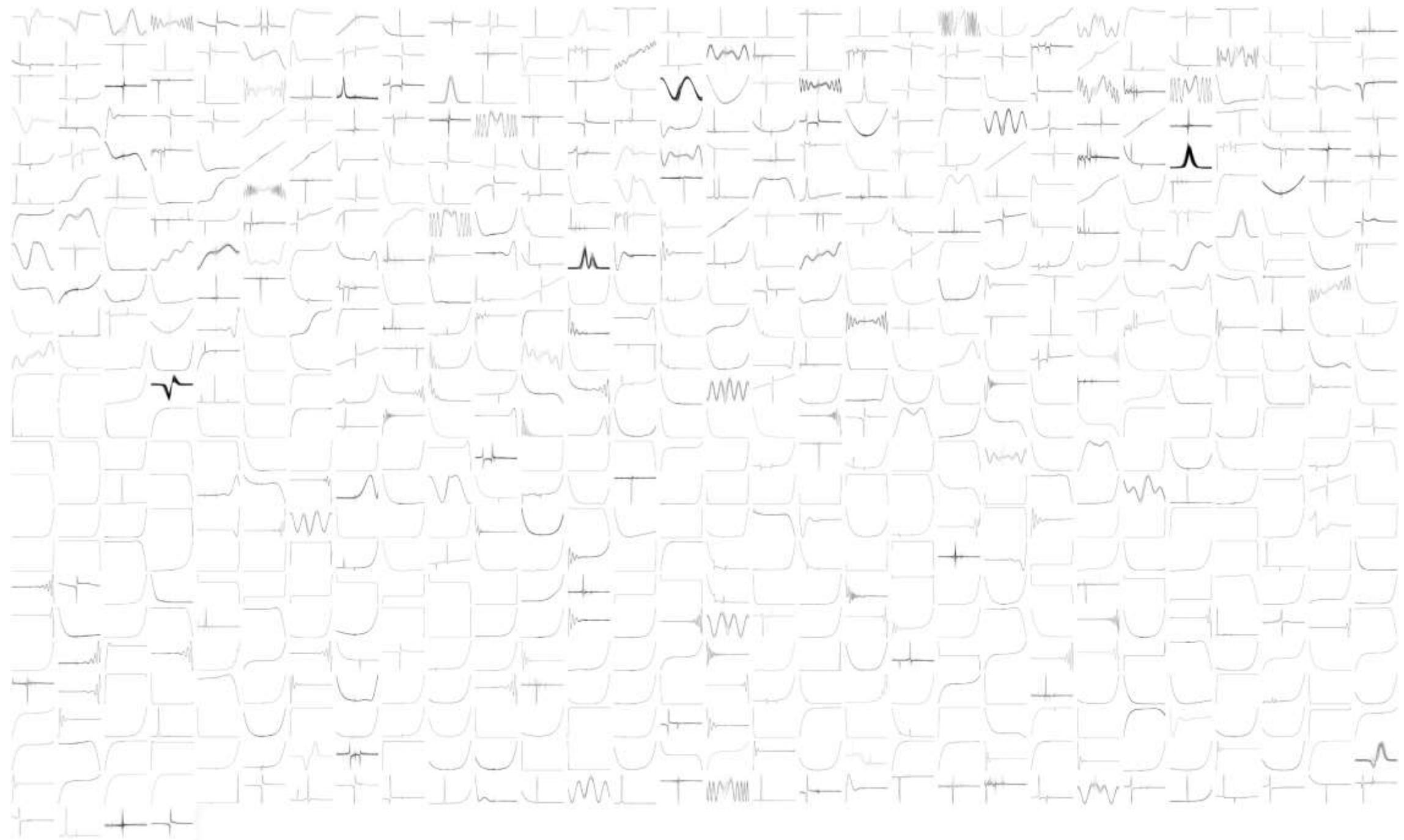
$$\frac{1}{\log(x) + xy^3}$$

How are symbolic regression models distributed in the search space?

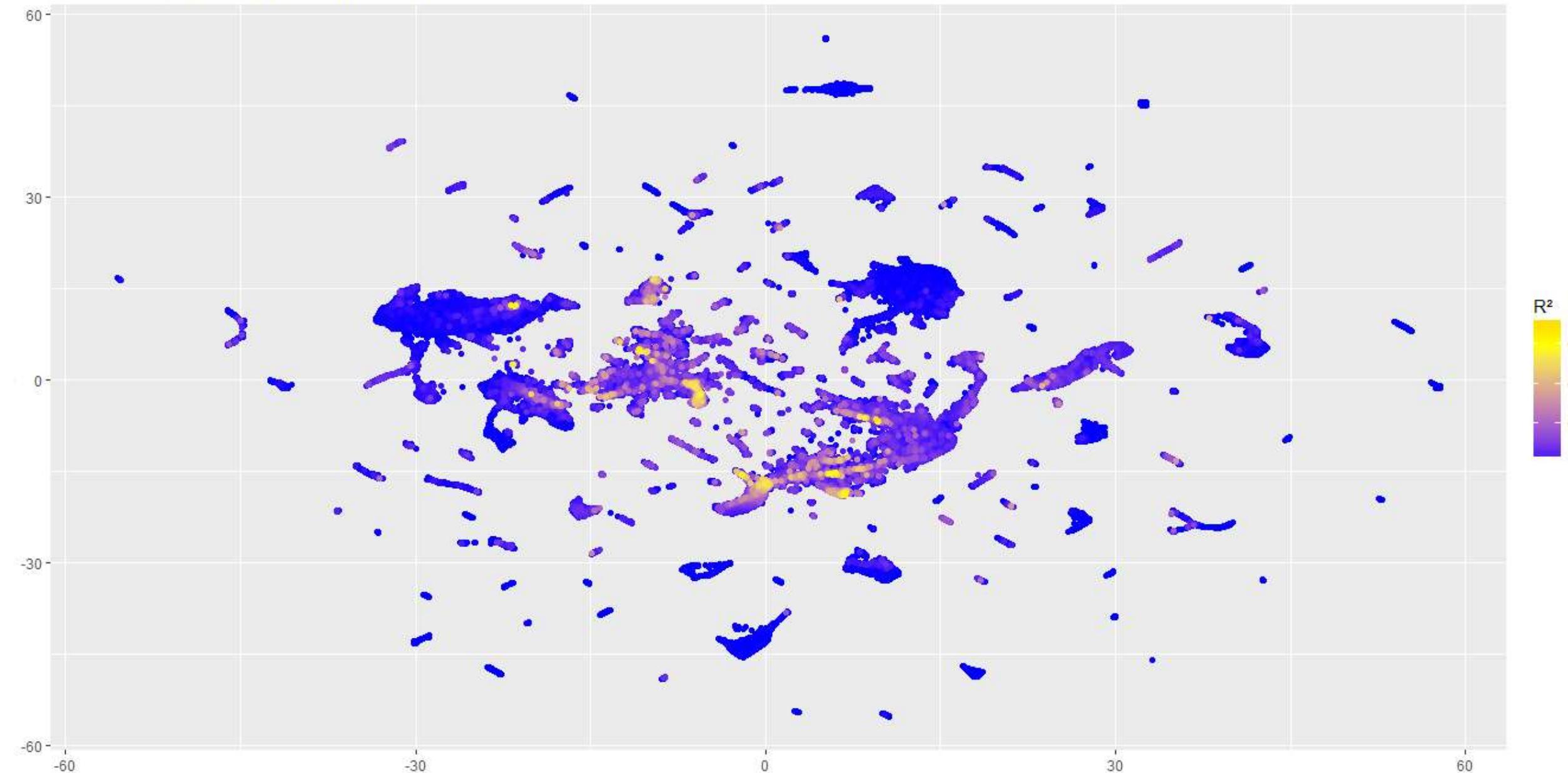
Which parts of the search space are visited by genetic programming?



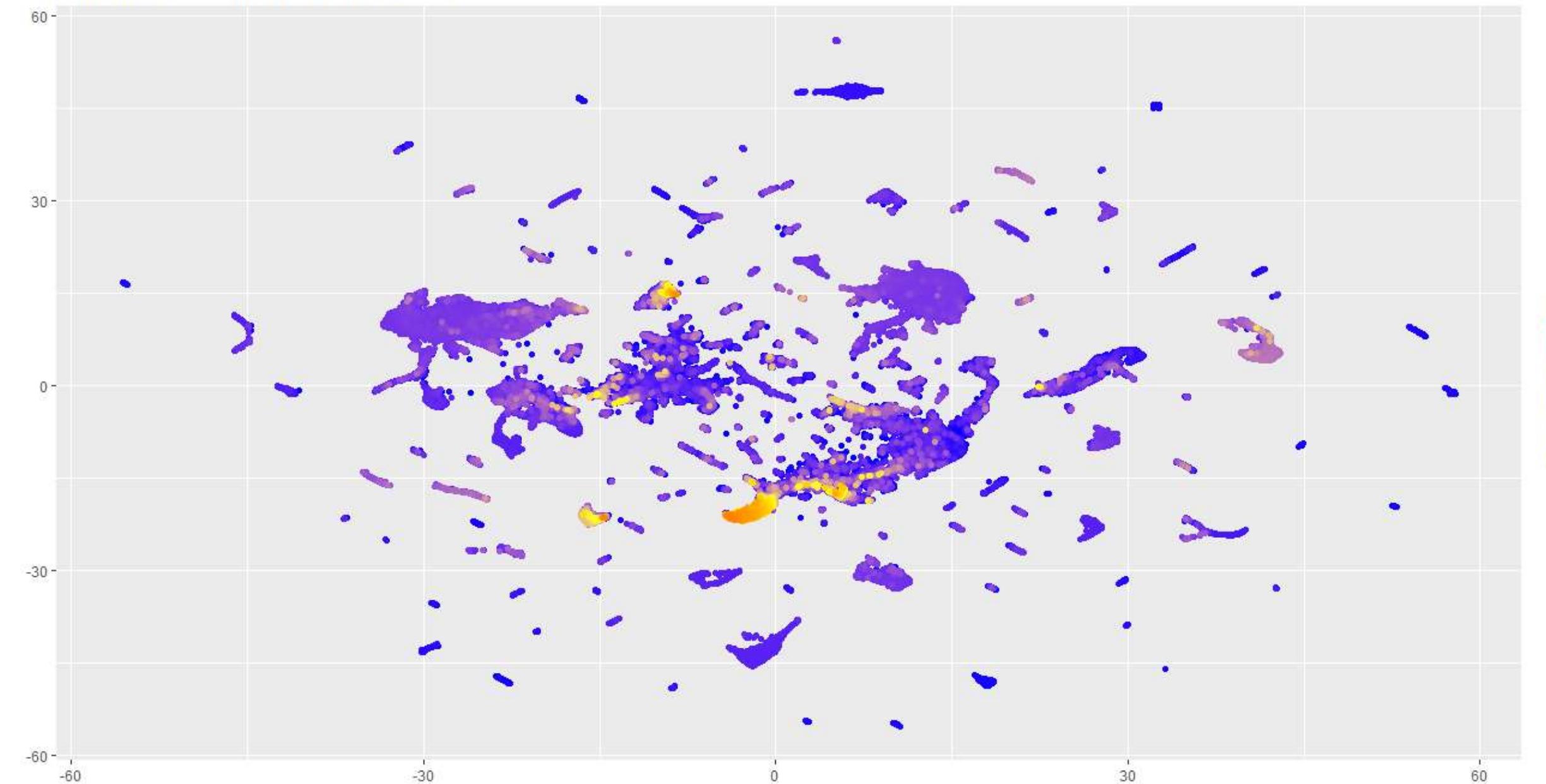




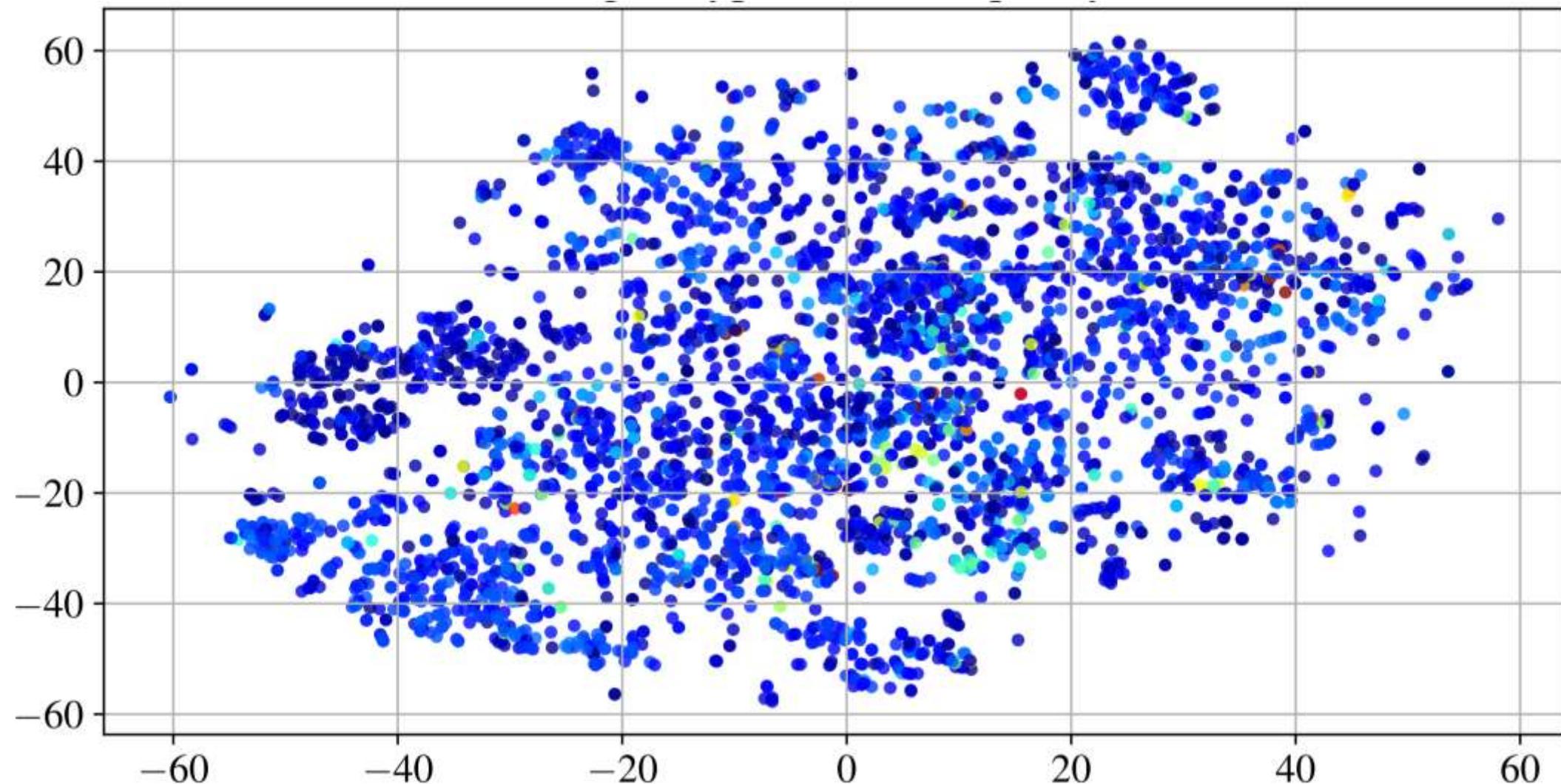
Phenotypic Embedding (R^2 with Keijzer4)



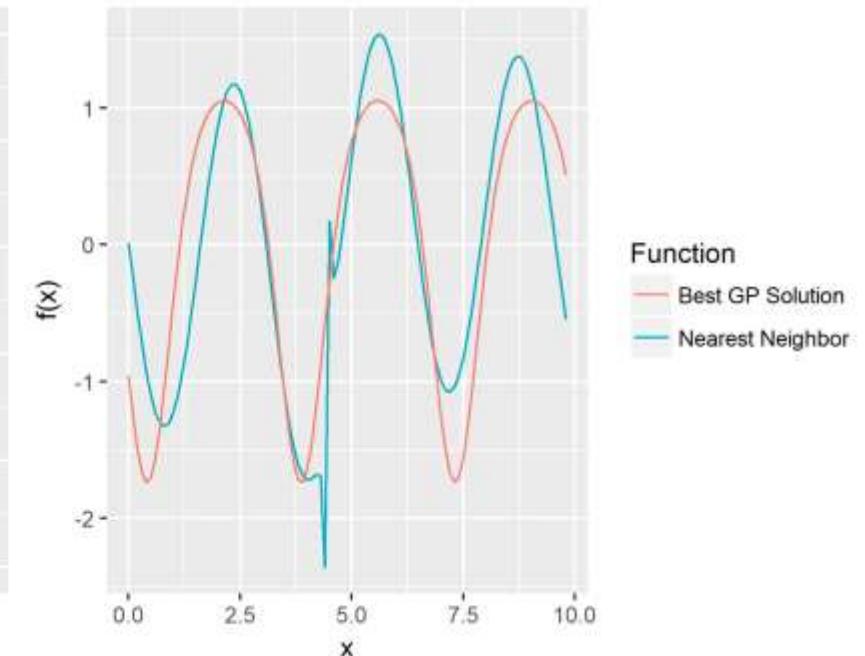
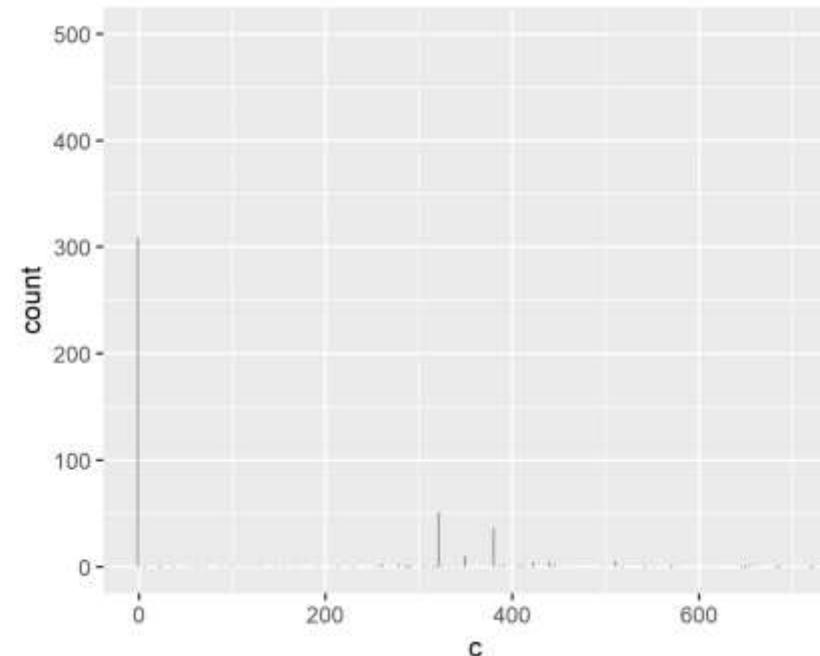
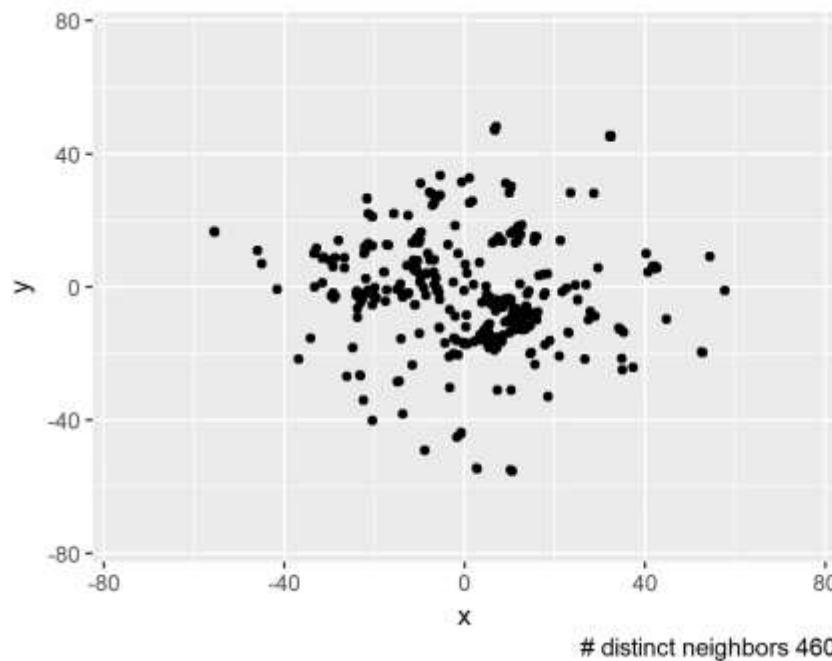
Phenotypic Embedding (R^2 with Pagie (1d))



Genotypic Embedding

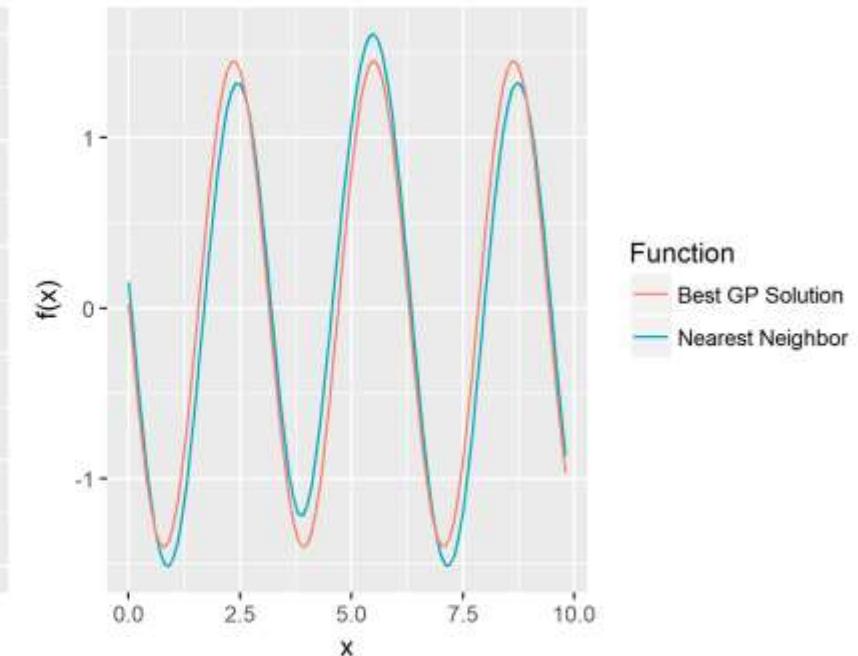
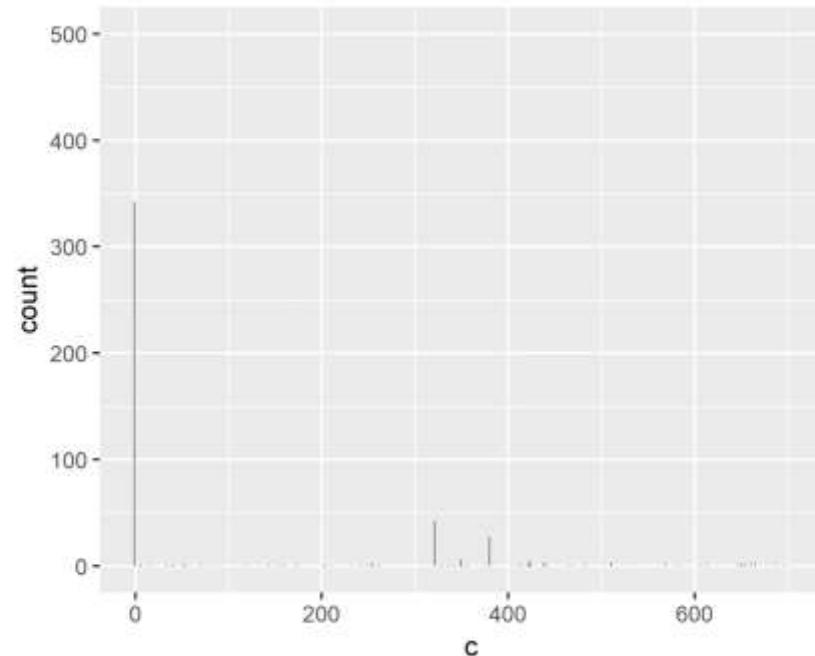
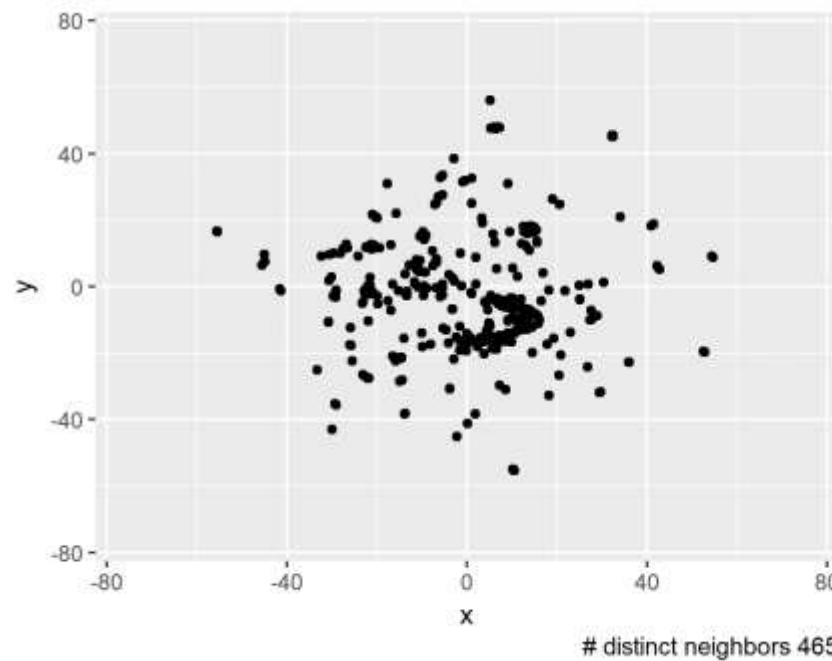


Mapping Populations



Standard GP, PopSize: 500 , Tournament: 7
https://youtu.be/Xtx_NKioazo

Mapping Populations



GP with offspring selection, PopSize: 500 , Random
<https://youtu.be/iYJEQXTPadE>

Discussion

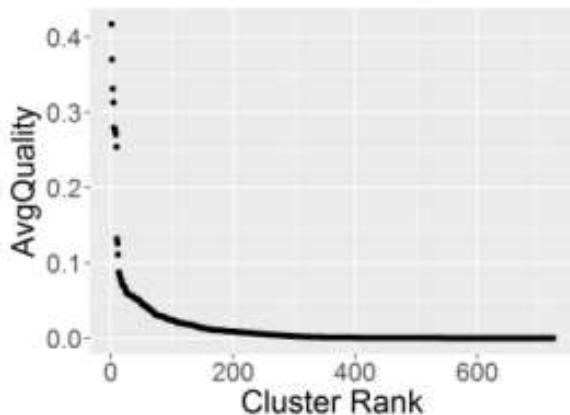
- What are the reasons for the low correlation between phenotypic and genotypic similarities?
 - Is it possible to apply the same approach to multi-variate models?
 - How should we include parameters?
 - How could we use the neighborhood graph in GP?
 - Are the insights transferable to real-world problems?
- ...



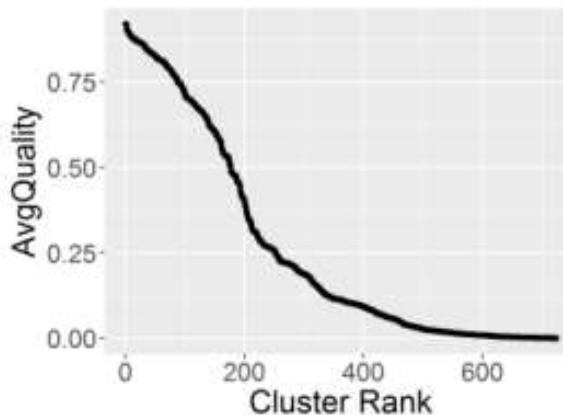
SymReg

**JOSEF RESSEL CENTER FOR
SYMBOLIC REGRESSION**

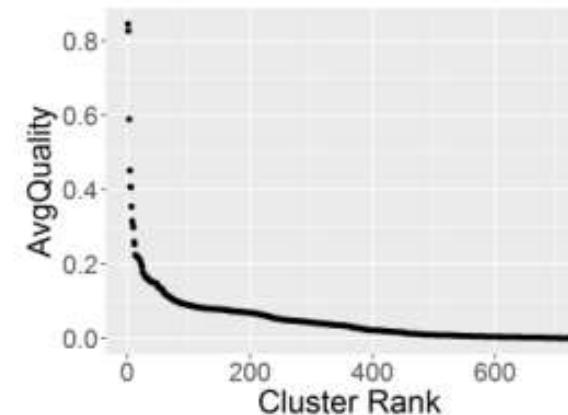




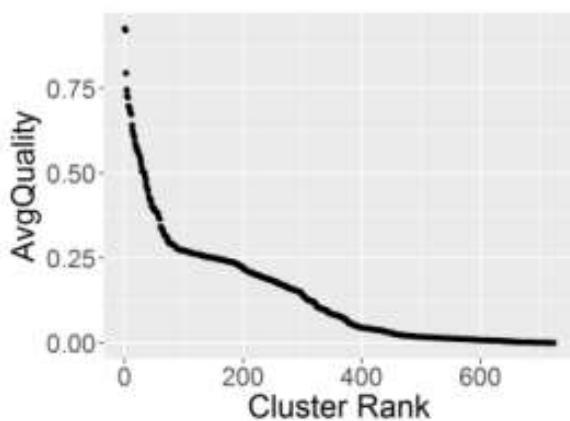
(a) Keijzer 4



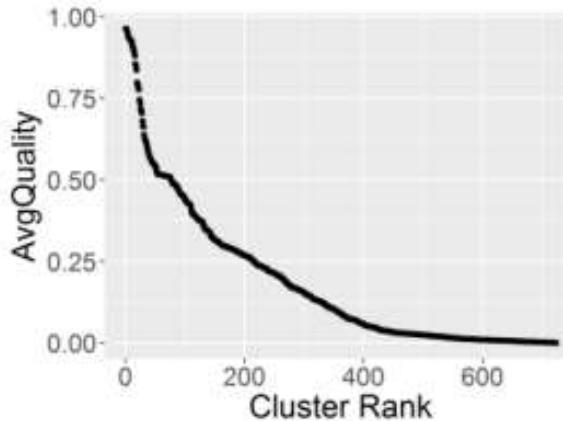
(b) Keijzer 9



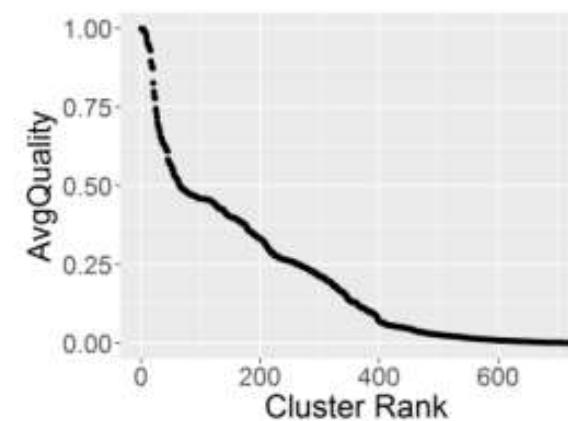
(c) Adapted Pagie



(d) Nguyen 5

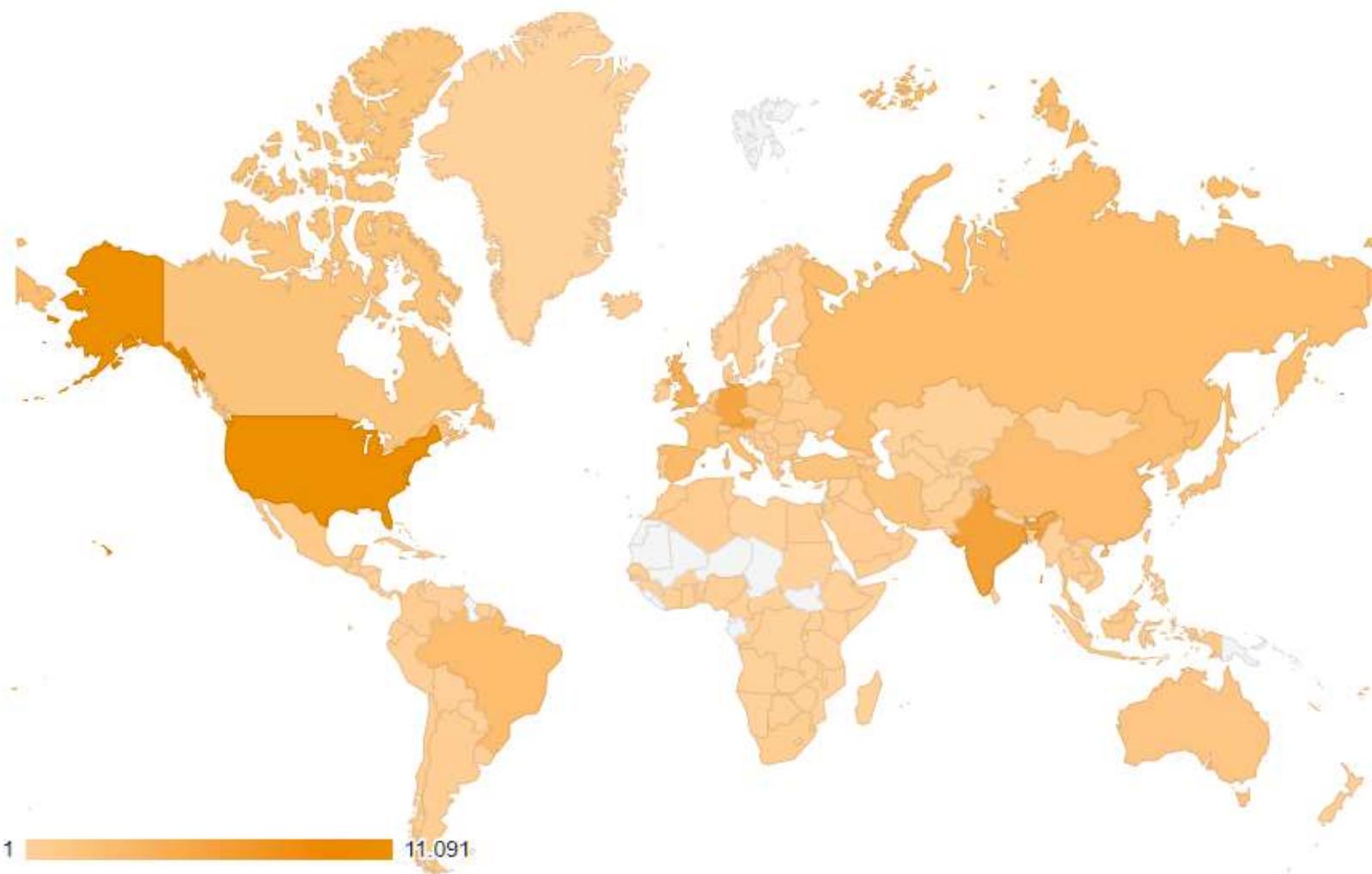


(e) Nguyen 6



(f) Nguyen 7

Fig. 5 Ranking of clusters by average R^2 values of expressions within all clusters. For each of the benchmark functions there are clusters which contain well-fitting expressions.



Land	Nutzer (dev.heuristiclab.com)
USA	11091
Austria	5766
India	5678
Germany	5096
UK	3743
Spain	2809
Russia	2645
Brazil	2607
China	2575
France	2099

Stand: Feb. 2018