



Published in final edited form as:

J Med Chem. 2011 October 13; 54(19): 6492–6500. doi:10.1021/jm200114f.

Novel Peptide-specific QSAR Analysis Applied to Collagen IV Peptides with Antiangiogenic Activity

Corban G. Rivera^{1,2,*}, Elena V. Rosca¹, Niranjana B. Pandey¹, Jacob E. Koskimaki¹, Joel S. Bader^{1,2}, and Aleksander S. Popel¹

¹Department of Biomedical Engineering, Johns Hopkins University, 720 Rutland Avenue, 613 Traylor Bldg, Baltimore, MD 21205

²High-Throughput Biology Center, Johns Hopkins School of Medicine, 733 N. Broadway Street Broadway Research Building, Room 359. Baltimore, MD 21205

Abstract

Angiogenesis is the growth of new blood vessels from existing vasculature. Excessive vascularization is associated with a number of diseases including cancer. Anti-angiogenic therapies have the potential to stunt cancer progression. Peptides derived from type IV collagen are potent inhibitors of angiogenesis. We wanted to gain a better understanding of collagen IV structure-activity relationships using a ligand-based approach. We developed novel peptide-specific QSAR models to study the activity of the peptides in endothelial cell proliferation, migration, and adhesion inhibition assays. We found that the models produced quantitatively accurate predictions of activity and provided insight into collagen IV derived peptide structure-activity relationships.

Background

Excessive vascularization is a hallmark of many diseases including cancer, rheumatoid arthritis, diabetic nephropathy, pathologic obesity, age-related macular degeneration, and asthma. Compounds that inhibit angiogenesis represent potential therapeutics for many diseases. Judah Folkman performed pioneering research in the field of angiogenesis;¹ his work led to the identification of a number of polypeptides with anti-angiogenic activity.² One of polypeptides called endostatin was derived from the noncollagenous (NC1) domain of collagen XVIII.³ Work led by Raghu Kalluri resulted in the development of small antiangiogenic peptides from the NC1 domain of collagen IV including canstatin,⁴ arrestin,⁵ and tumstatin.⁶ These collagen IV derived fragments were reviewed in the context of other angiogenesis modulating compounds.^{7–9} Based on these parent compounds, work in our laboratory identified more than 100 similar peptide sequences from diverse parent proteins throughout the proteome.¹⁰ The set of parent proteins included collagen IV, CXC chemokines, type I thrombospondin domain (TSP-1)-containing proteins, serpins, somatotropins, and tissue inhibitors of metalloproteinases (TIMPs). Work carried out in our group experimentally validated *in vitro* inhibition of endothelial cell (EC) proliferation and migration by peptides derived from type IV collagens,¹¹ thrombospondin domain-containing

*Corresponding author. Telephone: 410-955-1787, Fax 410-614-8796, cgrivera@jhu.edu.

Authors' contributions

CGR designed the method, performed the analysis, and wrote the paper. EVR, JEK and NBP performed the *in vitro* experiments. JSB and ASP motivated the problem, provided guidance for the analysis and manuscript. All co-authors edited the paper.

Competing Interests

The authors' declare no competing interests.

proteins,^{12, 13} and CXC chemokines.¹⁴ These studies showed that a large fraction of the peptides have antiangiogenic potential. Subsequently, our laboratory tested some of these peptides *in vivo* using mouse xenograft models of breast and lung cancer,^{15, 16} and ocular models.¹⁷ The peptides derived from type IV collagen are attractive targets because of their efficacy against multiple angiogenic properties (i.e. endothelial cell proliferation, migration, and adhesion).¹⁸

A better understanding of the structure-activity relationship of type IV collagen peptides could help us better understand the mechanism of action and produce more active peptides. For many of these peptides, the receptor had not been elucidated. When the receptor is unknown, ligand-based modeling approaches must be used. Examples of ligand-based design methods include pharamcophore modeling^{19–22} and quantitative structure-activity relationship (QSAR)^{23–26} analysis. These methods correlate diverse aspects of molecular structure and flexibility with a quantitative measure of activity. Some work has been done on developing peptide-specific feature sets for QSAR.^{27, 28} Others make use of position weight matrices to describe a family of peptides.²⁹ Many of these methods require solving NP-hard³⁰ problems. That means a polynomial time algorithm is not known for solving these problems. For large datasets, these methods must resort to using inexact approaches and heuristics.

To continue developing the type IV collagen-derived peptides, we aimed to (i) develop techniques for computationally efficient, peptide-specific, QSAR analysis, (ii) enable predictions of peptide activity, and (iii) gain a better understanding of the structure-activity relationship of collagen IV derived peptides. In this work, we described several novel peptide-specific QSAR methods that helped us address these aims. We formulated the models using convex optimization in a way that could be solved quickly to global optimality. We used experimentally-determined activity data from collagen IV peptides to develop individual models for endothelial cell proliferation, migration, and adhesion. We validated the QSAR models by making activity predictions and performing experiments for an external set of peptides. The activity of the external set of peptides was verified by endothelial cell proliferation, migration, adhesion, and tube formation assays.

Results

Peptide activity *in vitro* using EC proliferation, migration, and adhesion assays

This study is based on a library of 23 collagen IV derived peptides. The founding peptide **0** (SP2000)¹⁰ was found as a homolog of tumstatin⁶ in the human proteome. These peptides consisted of a series of truncations and selected amino acid substitutions designed to improve translational potential. In Table 1 we present the activity of the 23 (21 training + 2 external verification) peptides in endothelial cell proliferation (at 100 μ M), migration (at 50 μ M), and adhesion (at 100 μ M). Peptide concentrations were chosen to provide diversity in activity measurements. All experiments were performed in duplicate and the result of each experiment was the average of three replicates on the same plate. Activity measurements are given as a percentage of the vehicle control.

Modeling overview

In Figure 1, we outline the peptide modeling procedure. The methods are based on data that associates peptide features with a quantitative activity score (e.g., endothelial cell (EC) proliferation inhibition activity). Peptides are converted into unique sparse vector of features. For example, Figure 2 shows the vectorization of the short peptide LRRFSTMPFMF. In the simplest methodology that we consider, each feature uniquely identifies an amino acid at a single position. We use convex optimization to select features

that differentiate highly active and inactive peptides. We formulate the convex optimization objective in a way that can be solved quickly to global optimality.

Peptide-specific QSAR Method Comparison

We developed four approaches to model the data in Table 1 and learn about the structure-activity relationship of type IV collagen peptides. The approaches were based on the least absolute shrinkage and selection operator (Lasso).³¹ The approaches differed in the features that they consider and the weight assigned to training examples. The specific details of these approaches can be found in the Materials and Methods section.

In Table 2, we compared four methods for their ability to predict peptide efficacy. We compared each of these methods to a naive featureless method that always predicted the average activity from the training set. The methods were evaluated on three datasets that measured the ability of peptides to inhibit endothelial cell proliferation (A), migration (B), and adhesion (C). To compare these approaches, we took a leave-one-out cross validation (LOOCV) approach. The concept of LOOCV is that we use all but a single peptide to train the model. We then use that model to predict the efficacy of the single peptide, which was left out. This allowed us to compute the error between predicted and observed activity measurements. To determine which methods were statistically superior to others, we conducted t-tests for all pairs of methods based on their squared test errors. Significantly low test errors indicate better performance. The table gives the *p*-value associated two-tailed paired t-test. At the 0.05 level, all of the models had lower error than the naive featureless method. Also, the non-linear Lasso method had significantly less error than the Lasso method. These results held over all three datasets. Based on these results, the rest of the study was performed using non-linear Lasso. In Figure 3, we show the observed and leave-one-out predictions for each method for all peptides in the dataset in endothelial cell proliferation, migration, and adhesion assays. The figure illustrates that no single method had the least error in all trials, and that the predictive performance is good even in cases where percent inhibition is negative as seen in the migration and adhesion datasets.

QSAR analysis for type IV collagen derived peptides

In the previous sections we make extensive use of leave-one-out cross validations to estimate generalization error. We concluded from these analyses that non-linear Lasso had statistically lower generalization error than Lasso. Low generalization error is an indication that the features used in the models may be useful for understanding the structure-activity relationship of type IV collagen peptides.

In this section and unlike the previous sections, we train models for endothelial cell proliferation, migration, and adhesion based on all of the data in Table 1 except for the external validation set consisting of **27** and **35**. The models are structured such that important features receive high weight. The model features (first column) and weights (second column) are given in decreasing order in Table 3. The features are indicated for each row by the change in sequence from the preceding row. The weights were determined using the non-linear Lasso method (as described in Materials and Methods). We analyse these features for QSAR analysis. This approach gives us a way of indirectly identifying putative pharmacophores for the collagen-IV derived peptides.

When multiple amino acids are viable options in a position, they are shown in decreasing order of importance. In the migration model (Table 3, C) in the 18th position, L- α -amino-n-butyric acid (indicated by X) is preferred with a weight 0.018 over alanine with a weight of 0.016. The proliferation model (Table 3, A) makes it clear that there are important regions on the N-terminus (LRRF) and the C-terminus (NINNVXN). In the adhesion model (Table

3, B), the highly weighted asterisks in the 20th position indicates that truncation of the phenylalanine may improve the anti-adhesion activity of the peptide. Like the proliferation model, the regions on the N-terminus (LRRF) and C-terminus (NINNVX) are selected. Unlike the proliferation model, the L- α -amino-n-butyric acid in the 12th position is one of the most important features for anti-adhesion activity. The migration model (Table 3, C) highlights the C-terminal (ANINNVXN) as a useful indicator of anti-migration activity; however for full anti-migration activity the LRRF sequence is also required. From all three models we found that both the C-terminal sequence LRRF and the N-terminal sequence XNINNVXN are required for full activity.

Structural association

We examined the structure of peptide **0** as it exists in the native type IV collagen NC1 domain (pdb:1T60). In Figure 4, we show the conformation of the peptide in the native protein. By computing the solvent accessible surfaces of the protein, we found two exposed regions corresponding to the N-terminal (LRR) and C-terminal (INN). These regions correlate with the peptide motifs needed for anti-angiogenic activity.

Experimental model validation

Two peptides, **27** and **35**, were held out as an external validation set. Models for proliferation, migration, and adhesion were trained using all other peptides from Table 1. Based on these models, peptides **27** and **35** were predicted to have similar activity. They were predicted to have 54.15, 93.35, and 97.54 percent proliferation, migration, and adhesion inhibition, respectively. Based on the experimentally determined activities given in Table 1 and predicted activities, R^2 values on the external validation set were 0.84, 0.85, and 0.99 for the proliferation, migration, and adhesion models, respectively.³² From the R^2 values on the external validation set, we could conclude that the models were predictive for anti-angiogenesis phenotypes. In Figure 5, endothelial cell tube formation assays at 100 μ M confirmed the potency of peptides **27** (Figure 5, C) and **35** (Figure 5, D), relative to a vehicle control (Figure 5, A) and a weaker peptide **8** (SP2008) (Figure 5, B).

Discussion and Conclusions

Type IV collagens are basement membrane proteins that are essential for binding cells to the extracellular matrix.³³ Type IV collagen derived peptides have proven to be effective inhibitors of angiogenesis.³⁴ Using the models trained using the data from Table 1, we found a pair of regions namely LRRF at the C-terminus and XNINNVXN at the N-terminus are needed for full activity. This pair of important regions indicates that secondary structure or multiple binding sites may be important for the endothelial cell proliferation, migration, and adhesion inhibition activity of type IV collagen derived peptides. These results are consistent with a previous study on the tumstatin peptide by Eikesdal *et al.*³⁵ They found that the mutations to the NINN region resulted in a significant change in EC proliferation inhibition. These results also indicate that truncations to the 20-mer peptide with the exception of the phenylalanine in the 20th position would be detrimental to the activity of the collagen IV derived peptides.

In this article, we describe four novel peptide-specific QSAR approaches. We compared these approaches by testing their ability to predict the outcome of *in vitro* experiments. The comparison indicated that one approach called non-linear Lasso had statistically lower generalization error than Lasso (Table 2). We showed the individual predictions made by this approach in Figure 3. We found that the predictions made using the all four approaches were statistically significant compared to a method based on naive predictions. These results gave us confidence in the utility of the peptide-specific QSAR models. We analyzed the

features of these models to learn about the structure-activity relationship of collagen IV derived peptides. By analysing the structure of the collagen IV NC1 domain, we found that the solvent accessible regions of the peptide in the parent protein correlated with the motifs needed for anti-angiogenic activity.

Materials and methods

Peptide dataset

All peptides were synthesized by New England Peptide with at least 95% purity evaluated using both HPLC and MALDI by the manufacturer. Table 1 gives the compound structures in terms of the one letter amino acid codes. Truncated amino acids are indicated by asterisks. The error in the activity measurements was based on two biological replicates each derived from the mean of three technical replicates. The data are shown as percent inhibition relative to a vehicle control. A single dose was selected for each dataset that produced a diverse set of activities for the candidate peptides. Proliferation and adhesion measurements were taken at a peptide concentration of 100 μ M, while migration measurements were taken with a compound dose of 50 μ M.

Cell culture

Human umbilical vein endothelial cells (HUVEC) were purchased from Lonza and were grown under the manufacturer's recommendation using Endothelial Basal Media (EBM-2) supplemented with the Bullet Kit (EGM-2, Lonza). Cells of passages 2–7 were used for experiments. Cells were grown at 37°C in a humidified incubator with 5% CO₂.

Proliferation assays

Colorimetric WST-1 reagent (Roche, IN) was used to perform the proliferation assays. HUVECs were plated in 96-well plates at a 2000 cell/well density. Peptides at 100 μ M in fully supplemented media were added to the adherent cells and incubated for 72 hours. WST-1 reagent was added in serum free media for four hours and the color intensity was measured at 450 nm with Victor-V plate reader (Perkin Elmer, MA).

Migration assay

The effect of the migration inhibition of the peptides on the cells was determined using electrical impedance measurements with a continuous and real time migration assay (RT-CIM, ACEA Biosciences, CA). The top compartment of the CIM plate was coated with fibronectin (20 μ g/ml) and 45,000 HUVEC/well were added either in the presence or absence of the peptide at 50 μ M. Fully supplemented media was added to the bottom compartment serving as chemoattractant. The migration of the cells is measured by the integrated sensors in the bottom side of the porous membrane which divides the two chambers. This technology allows for easy quantification of cell migration by monitoring the cell index (derived from the measured impedances).

Adhesion assay

The adhesion inhibitory potential of the peptides was also measured using RT-CIM technology. In this instance single compartment E-plates (ACEA, Biosciences, CA) were used, in which 25,000 HUVEC/well were plated in the presence or absence of the peptides at 100 μ M and the adhesion measured by the changes in the cell index amplitude for 3 hours.

Tube formation

Tube formation assay was performed by following the published protocol by Arnaoutva *et al.*³⁶ Briefly, 96 well plates were coated with Geltrex, Reduced Growth Factor Basement

Membrane Matrix (Invitrogen, CA) (50 μ l/well) and incubated at 37°C for 30 minutes to allow gelation to occur. HUVECs were added to the top of the gel at a density of 15,000 cells/well in the presence or absence of the peptide (100 μ M). The positive control included the same amount of solvation vehicle (i.e., DMSO) as the experimental condition. Cells were incubated at 37°C with 5% CO₂ overnight and pictures were captured with a CCD Sensicam camera mounted on a Nikon inverted microscope.

Peptide-specific QSAR approaches

We took as input a set of peptide sequences along with an experimentally measured efficacy for each peptide. The method returned a model which could be used to predict the efficacy of hypothetical peptides from the same class. The method worked by converting each peptide sequence into an input space of amino acids and positions. Those were the explanatory variables in the peptide-specific QSAR modeling framework. A weight for each feature was learned using non-negative Lasso regression³⁷ with the peptide efficacies as response variables. The scaling term for the L1-norm regularizer was determined using leave-one-out cross validation. Despite evaluating many features, the use of L1-norm regularization allowed the model to avoid over-fitting. The convex nature of the optimization problem allowed the method to quickly reach the globally optimal solution without a combinatorial search of input space. The software which was implemented in Matlab using CVX³⁸ is freely available upon request.

Lasso with an amino acid substitution matrix

Without loss of generality, we describe the method in terms of the 20 common amino acids. Given m peptides of length n , let p_{ij} be amino acid j in peptide i . Let r be a list of all 20 natural amino acids. Let \mathbf{S} be a 20 by 20 amino acid association matrix, in this study we use the PAM250 matrix,³⁹ such that $\mathbf{S}(a,b)$ gives the association between amino acid a and b . We use the PAM250 matrix as a principled approach to give weight to amino acids with similar biochemical properties. Let \mathbf{A} be an m by $20n$ matrix that encodes the amino acid sequences, such that

$$\mathbf{A}_{i,jk} = \mathbf{S}(p_{i,j}, r_k) \quad (1)$$

Let \mathbf{b} be a vector of length m representing the activity of each peptide. In this study, the quantitative measure of activity is given by percent endothelial cell proliferation, migration, or adhesion inhibition. Our goal is to learn values in the weight vector \mathbf{x} of length $20n$. The values in the weight vector \mathbf{x} correspond to the relative importance of the features considered in the model. Using this formulation, we solve the standard Lasso objective subject to $\mathbf{x} \geq 0$. Lasso is composed of the least-squares objective regularized by the L1 norm of the weight vector. The parameter λ influences the sparsity of the weight vector \mathbf{x}

$$\min \cdot \|\mathbf{Ax} - \mathbf{b}\|_2 + \lambda \|\mathbf{x}\|_1 \quad (2)$$

Non-linear Lasso

In the previous section, we described the linear version of Lasso using only the input space described in \mathbf{A} . As an alternate approach, we expand on the input space given in the previous approach to a feature space consisting of pairs of features. Let \mathbf{A}' be an m by $(20n)^2$ matrix. Although the number of features is large, we use sparse matrices to eliminate unused variables and reduce the problem size. We make use of aggressive regularization to avoid

over-fitting. The Lagrange multiplier λ is selected automatically by leave-one-out cross validation. We use the objective from equation (2) except that we make use of \mathbf{A}' and the \mathbf{x} vector is of length $(20n)^2$.

Locally-weighted methods

We extend both linear and non-linear Lasso to construct locally-weighted variants of both methods. The idea is that we will weight training examples in \mathbf{A} by their proximity to the vectorized peptide \mathbf{y} to be predicted. The intuition is that we prefer to make smaller training errors for points close to the test point \mathbf{y} . The weight \mathbf{w} assigned to each training example in \mathbf{A} is given in equation (3).

$$\mathbf{w}_j = \exp \left[-\|\mathbf{y} - \mathbf{A}_j\|_2 \right] \quad (3)$$

The weighted objective for the linear version of Lasso is given in equation (4).

$$\min \cdot \sum_i w_i (\mathbf{A}_i \mathbf{x} - b_i)^2 + \lambda \|\mathbf{x}\|_1 \quad (4)$$

Statistical significance and cross validation

To evaluate the quality of the predictions given by the peptide-specific QSAR approaches, we perform leave-one-out cross validation. For each of the m peptide examples, we split the examples into a test set containing the i^{th} peptide and a training set containing all other peptides. We use the training set of peptides to obtain the weight vector \mathbf{x} . Let \mathbf{p}_i be the vector of length $20n$ that encodes the i^{th} peptide. The predicted activity q_i for the i^{th} peptide is given by

$$\mathbf{q}_i = \mathbf{p}_i^T \mathbf{x} \quad (5)$$

The statistical significance of the predictions is determined by comparing the set of residuals generated using our model predictions with residuals generated using naive model predictions. We test the null hypothesis that the residuals between the observed and predicted values are equal to the residuals between the observed and naive model predictions (i.e., a model that always predicts the mean training efficacy). The alternative hypothesis is that the residuals between the observed and predicted values are less than the residuals between the observed and naive model predictions. We generate a p -value for each model using a one-sided paired t-test. We used R^2 as a metric of model performance on the external validation set.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y - \hat{y})^2}{\sum_{i=1}^n (y - \bar{y})^2} \quad (6)$$

In this metric, experimentally observed values y are compared with predicted values \hat{y} relative to the mean observed value from the training set \bar{y} .

Acknowledgments

The work was supported by NIH grants R01 CA138264, R01 HL101200, and U54RR020839.

Abbreviations

QSAR	quantitative structure activity relationship
NC1	non-collagenous domain
CXC	n-terminal cysteine-x- cysteine domain
TSP-1	thrombospondin 1 domain
TIMP	tissue inhibitor of metalloproteinases
EC	endothelial cell
NP-hard	non-deterministic polynomial-time hard
Lasso	least absolute shrinkage and selection operator
LOOCV	leave-one-out cross validation
HUVEC	human umbilical vein endothelial cells
L1-norm	one norm
PAM250	250% point accepted mutation matrix

References

1. Folkman J. Tumor angiogenesis: therapeutic implications. *N Engl J Med.* 1971; 285:1182–1186. [PubMed: 4938153]
2. Folkman J. Angiogenesis: an organizing principle for drug discovery? *Nat Rev Drug Discov.* 2007; 6:273–286. [PubMed: 17396134]
3. O'Reilly MS, Boehm T, Shing Y, Fukai N, Vasios G, Lane WS, Flynn E, Birkhead JR, Olsen BR, Folkman J. Endostatin: an endogenous inhibitor of angiogenesis and tumor growth. *Cell.* 1997; 88:277–285. [PubMed: 9008168]
4. Kamphaus GD, Colorado PC, Panka DJ, Hopfer H, Ramchandran R, Torre A, Maeshima Y, Mier JW, Sukhatme VP, Kalluri R. Canstatin, a novel matrix-derived inhibitor of angiogenesis and tumor growth. *Journal of Biological Chemistry.* 2000; 275:1209. [PubMed: 10625665]
5. Nyberg P, Xie L, Sugimoto H, Colorado P, Sund M, Holthaus K, Sudhakar A, Salo T, Kalluri R. Characterization of the anti-angiogenic properties of arresten, an $[\alpha] 1 [\beta] 1$ integrin-dependent collagen-derived tumor suppressor. *Experimental cell research.* 2008; 314:3292–3305. [PubMed: 18775695]
6. Maeshima Y, Sudhakar A, Lively JC, Ueki K, Kharbanda S, Kahn CR, Sonenberg N, Hynes RO, Kalluri R. Tumstatin, an endothelial cell-specific inhibitor of protein synthesis. *Science.* 2002; 295:140. [PubMed: 11778052]
7. Kalluri R. Basement membranes: structure, assembly and role in tumour angiogenesis. *Nature Reviews Cancer.* 2003; 3:422–433.
8. Mundel TM, Kalluri R. Type IV collagen-derived angiogenesis inhibitors. *Microvascular research.* 2007; 74:85–89. [PubMed: 17602710]
9. Nyberg P, Xie L, Kalluri R. Endogenous inhibitors of angiogenesis. *Cancer research.* 2005; 65:3967. [PubMed: 15899784]
10. Karagiannis ED, Popel AS. A systematic methodology for proteome-wide identification of peptides inhibiting the proliferation and migration of endothelial cells. *Proc Natl Acad Sci U S A.* 2008; 105:13775–13780. [PubMed: 18780781]

11. Karagiannis ED, Popel AS. A theoretical model of type I collagen proteolysis by matrix metalloproteinase (MMP) 2 and membrane type 1 MMP in the presence of tissue inhibitor of metalloproteinase 2. *J Biol Chem.* 2004; 279:39105–39114. [PubMed: 15252025]
12. Karagiannis ED, Popel AS. Anti-angiogenic peptides identified in thrombospondin type I domains. *Biochem Biophys Res Commun.* 2007; 359:63–69. [PubMed: 17531201]
13. Karagiannis ED, Popel AS. Peptides derived from type I thrombospondin repeat-containing proteins of the CCN family inhibit proliferation and migration of endothelial cells. *Int J Biochem Cell Biol.* 2007; 39:2314–2323. [PubMed: 17707681]
14. Karagiannis ED, Popel AS. Novel anti-angiogenic peptides derived from ELR-containing CXC chemokines. *J Cell Biochem.* 2008; 104:1356–1363. [PubMed: 18307172]
15. Koskimaki JE, Karagiannis ED, Rosca EV, Vesuna F, Winnard PT Jr, Raman V, Bhujwala ZM, Popel AS. Peptides derived from type IV collagen, CXC chemokines, and thrombospondin-1 domain-containing proteins inhibit neovascularization and suppress tumor growth in MDA-MB-231 breast cancer xenografts. *Neoplasia.* 2009; 11:1285–1291. [PubMed: 20019836]
16. Koskimaki JE, Karagiannis ED, Tang BC, Hammers H, Watkins DN, Pili R, Popel AS. Pentastatin-1, a collagen IV derived 20-mer peptide, suppresses tumor growth in a small cell lung cancer xenograft model. *BMC Cancer.* 2010; 10:29. [PubMed: 20122172]
17. Cano Mdel V, Karagiannis ED, Soliman M, Bakir B, Zhuang W, Popel AS, Gehlbach PL. A peptide derived from type 1 thrombospondin repeat-containing protein WISP-1 inhibits corneal and choroidal neovascularization. *Invest Ophthalmol Vis Sci.* 2009; 50:3840–3845. [PubMed: 19279315]
18. Rosca EV, Koskimaki JE, Rivera CG, Pandey NB, Tamiz AP, Popel AS. Anti-angiogenic peptides for cancer therapeutics. *Curr Pharm Biotechnol.* 2011; 12:1101–1116. [PubMed: 21470139]
19. Dixon SL, Smondirev AM, Knoll EH, Rao SN, Shaw DE, Friesner RA. PHASE: a new engine for pharmacophore perception, 3D QSAR model development, and 3D database screening: 1. Methodology and preliminary results. *J Comput Aided Mol Des.* 2006; 20:647–671. [PubMed: 17124629]
20. Schneidman-Duhovny D, Dror O, Inbar Y, Nussinov R, Wolfson HJ. PharmaGist: a webserver for ligand-based pharmacophore detection. *Nucleic Acids Res.* 2008; 36:W223–228. [PubMed: 18424800]
21. Güner, OF. Pharmacophore perception, development, and use in drug design. *Intl Univ Line*; 1999.
22. Mason JS, Morize I, Menard PR, Cheney DL, Hulme C, Labaudiniere RF. New 4-point pharmacophore method for molecular similarity and diversity applications: overview of the method and applications, including a novel approach to the design of combinatorial libraries containing privileged substructures. *Journal of medicinal chemistry.* 1999; 42:3251–3264. [PubMed: 10464012]
23. Blankley, C. Quantitative structure-activity relationships of drugs. Academic Press; New York: 1983. Introduction: A review of QSAR methodology.
24. Dudek AZ, Arodz T, Galvez J. Computational methods in developing quantitative structure-activity relationships (QSAR): a review. *Combinatorial Chemistry; High Throughput Screening.* 2006; 9:213–228.
25. Hansch, C.; Leo, A.; Hoekman, D. Exploring QSAR.:. Fundamentals and applications in chemistry and biology. An American Chemical Society Publication; 1995.
26. Kubiny H. Variable selection in QSAR studies. I. An evolutionary algorithm. *Quantitative Structure Activity Relationships.* 1994; 13:285–294.
27. Lin ZH, Long HX, Bo Z, Wang YQ, Wu YZ. New descriptors of amino acids and their application to peptide QSAR study. *Peptides.* 2008; 29:1798–1805. [PubMed: 18606203]
28. Zhou P, Chen X, Shang Z. Side-chain conformational space analysis (SCSA): a multi conformation-based QSAR approach for modeling and prediction of protein-peptide binding affinities. *J Comput Aided Mol Des.* 2009; 23:129–141. [PubMed: 18841329]
29. Doytchinova IA, Blythe MJ, Flower DR. Additive method for the prediction of protein-peptide binding affinity. Application to the MHC class I molecule HLA-A* 0201. *Journal of proteome research.* 2002; 1:263–272. [PubMed: 12645903]

30. Finn P, Halperin D, Kaviraki L, Latombe JC, Motwani R, Shelton C, Venkatasubramanian S. Geometric manipulation of flexible ligands. *Applied Computational Geometry Towards Geometric Engineering*. 1996:67–78.
31. Tibshirani R. Regression Shrinkage and Selection Via the Lasso. *Journal of the Royal Statistical Society, Series B*. 1994; 58:267–288.
32. Hawkins DM, Basak SC, Mills D. Assessing model fit by cross-validation. *Journal of chemical information and computer sciences*. 2003; 43:579–586. [PubMed: 12653524]
33. Khoshnoodi J, Pedchenko V, Hudson BG. Mammalian collagen IV. *Microscopy research and technique*. 2008; 71:357–370. [PubMed: 18219669]
34. Maeshima Y, Manfredi M, Reimer C, Holthaus KA, Hopfer H, Chandamuri BR, Kharbanda S, Kalluri R. Identification of the anti-angiogenic site within vascular basement membrane-derived tumstatin. *J Biol Chem*. 2001; 276:15240–15248. [PubMed: 11278365]
35. Eikesdal HP, Sugimoto H, Birrane G, Maeshima Y, Cooke VG, Kieran M, Kalluri R. Identification of amino acids essential for the antiangiogenic activity of tumstatin and its use in combination antitumor activity. *Proceedings of the National Academy of Sciences*. 2008; 105:15040.
36. Arnaoutova I, George J, Kleinman HK, Benton G. The endothelial cell tube formation assay on basement membrane turns 20: state of the science and the art. *Angiogenesis*. 2009; 12:267–274. [PubMed: 19399631]
37. Tibshirani R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (Methodological)*. 1996; 58:267–288.
38. Grant M, Boyd S. Graph implementations for nonsmooth convex programs. *Recent advances in learning and control*. 2008:95–110.
39. Jones DT, Taylor WR, Thornton JM. The rapid generation of mutation data matrices from protein sequences. *Computer applications in the biosciences: CABIOS*. 1992; 8:275. [PubMed: 1633570]

Peptide Vectorization with PAM250 augmentation

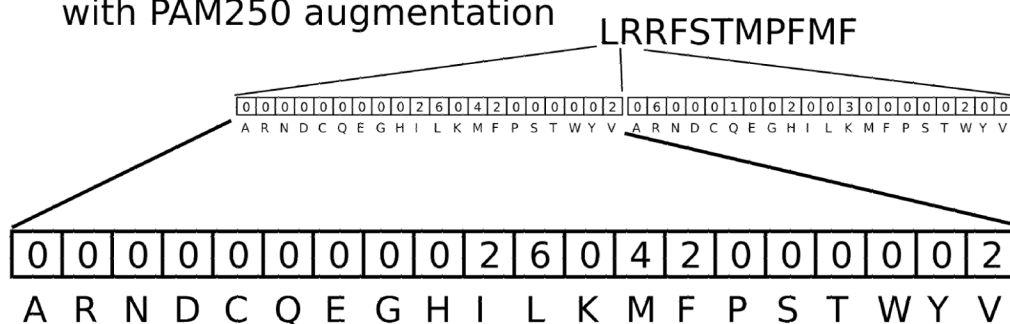


Figure 2. Peptide vectorization

Each peptide is converted into a sparse vector which uniquely maps specific amino acids to positions in the peptide. The mapping is augmented by the PAM250 amino acid substitution matrix. PAM matrices are based on the empirical mutation rate of amino acids in evolutionarily related proteins. For example, the figure shows the vectorization of the peptide LRRFSTMPFMF. The first amino acid leucine (L) can mutate to isoleucine (I), methionine (M), phenylalanine (F), and valine (V) at a rates greater than expected by chance. The weights assigned to these amino acids are given by log odds ratio in the PAM250 matrix. All other amino acids mutate from leucine at a lower rate than expected by chance. As a result, their value is set to zero. The PAM matrix gave us a principled way to associate common amino acids based on their chemical and structural properties.

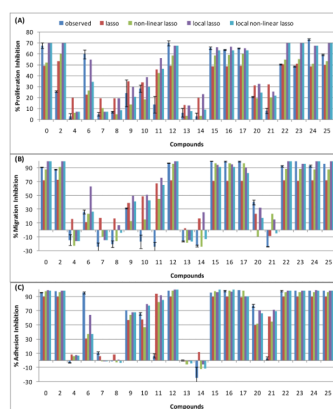


Figure 3. Quantitative predictions of peptide activity using non-linear Lasso

The observed and predicted activity of the 21 training peptides screened in endothelial cell (A) proliferation, (B) migration, and (C) adhesion assays. Compounds are given in Table 1. Predictions are made using LOOCV to assess the generalization error of the method. Predictions are shown for the four methods described in this Materials and Methods. The results imply an average error of between 14-20% depending on the assay.

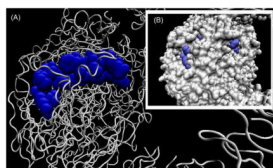


Figure 4. Solvent accessible surfaces of the peptide 0 in non-collagenous (NC1) domain of collagen IV

(A) the location of the peptide 0 in the NC1 domain of collagen IV. (B) the solvent exposed surfaces of peptide 0. The regions at the N-terminus and C-terminus are solvent accessible.

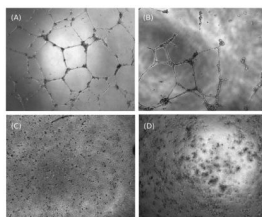


Figure 5. Endothelial tube formation assay

Endothelial cell tube formation assays are useful indicators of angiogenesis potential. (A) Tube formation for the positive control (vehicle control). HUVECs form robust tube structures (B) Endothelial cell tube formation with the addition of 100 μ M of **8**. The figure shows only partial inhibition of tube structures. (C) 100 μ M of **27** completely inhibits the formation of tube structures. (D) 100 μ M of **35** completely inhibits the formation of tube structures.

Table 1

The compound database

A dataset of 23 collagen IV derived compounds tested for endothelial cell proliferation (at 100 mg/ml), migration (at 50 mg/ml), and adhesion (at 100 mg/ml) inhibition. The table gives the mean % inhibition for each assay and the standard error of the mean (SEM). The screening of all compounds was done with n=2 and normalized to a vehicle control. We use the single letter X to represent L- α -amino-n-butyric acid (Abu). Peptides **27** and **35** were held out as an external validation set.

Ref.	Compound	Structure	Proliferation	Migration	Adhesion
0	SP2000	LRRFSTMPFMFCNINNVCFN	67.55 \pm 2.58	91.02 \pm 0.63	96.25 \pm 0.34
2	SP2002	LRRFSTMPFMFGNINNVGNF	25.71 \pm 0.78	87.90 \pm 0.70	97.80 \pm 0.15
4	SP2004	LRRFSTMPFMF*****	2.97 \pm 2.62	-14.20 \pm 7.85	-2.60 \pm 0.32
6	SP2006	LRRFSTMPFMFXNINNV****	59.70 \pm 4.10	26.55 \pm 3.15	95.62 \pm 1.22
7	SP2007	LRRFSTMPFMFX*****	4.85 \pm 1.62	-24.38 \pm 5.87	11.04 \pm 1.81
8	SP2008	LRRFSTMP*****	7.02 \pm 0.69	-20.09 \pm 4.67	-0.14 \pm 0.00
9	SP2009	*****NINNVXNF	24.25 \pm 12.25	31.65 \pm 0.85	70.78 \pm 0.05
10	SP2010	*****FMFXNINNVXNF	28.30 \pm 3.39	-16.60 \pm 9.83	66.18 \pm 1.83
11	SP2011	***STMFMFXNINNVXNF	13.83 \pm 7.64	-24.46 \pm 6.89	6.82 \pm 2.87
12	SP2012	LRRFSTMPFMFXNINNVXNF	69.72 \pm 2.25	96.97 \pm 0.19	98.46 \pm 0.32
13	SP2013	LNRFFSTMP*****	6.35 \pm 4.17	-16.05 \pm 0.69	0.32 \pm 0.86
14	SP2014	LRRFSTNLPFNLF*****	2.95 \pm 3.01	-22.77 \pm 1.28	-24.59 \pm 15.37
15	SP2015	LRRFSTMPAMFXNINNVXNF	65.40 \pm 1.20	99.75 \pm 0.25	95.60 \pm 0.31
16	SP2016	LRRFSTMPFAFXNINNVXNF	63.60 \pm 0.35	99.60 \pm 0.40	98.33 \pm 1.14
17	SP2017	LRRFSTMPFMA XNINNVXNF	65.10 \pm 0.57	99.60 \pm 0.40	99.27 \pm 0.21
20	SP2020	*****FXNINNVXN*	20.55 \pm 0.46	40.35 \pm 3.45	77.51 \pm 2.36
21	SP2021	*****FXNIN*****	8.17 \pm 2.27	-23.68 \pm 0.59	3.45 \pm 2.25
22	SP2022	LRRFSTMPFMFSNINNVSNF	50.47 \pm 0.66	92.56 \pm 0.97	97.95 \pm 1.27
23	SP2023	LRRFSTMPFMAINNVANF	48.61 \pm 0.53	99.59 \pm 0.13	98.19 \pm 0.03
24	SP2024	LRRFSTMPFMFINNVVINF	73.20 \pm 0.78	92.92 \pm 1.78	98.58 \pm 0.12
25	SP2025	LRRFSTMPFMFTNINNVTNF	59.10 \pm 0.85	96.08 \pm 0.15	98.69 \pm 0.08
27	SP2027	LRRFSTMPFMFVNINNVVNF	62.20 \pm 0.64	98.78 \pm 0.14	95.12 \pm 0.62
35	SP2035	LRRFSTMPFAFINNVVINF	46.58 \pm 4.23	69.33 \pm 0.14	94.80 \pm 0.93

Table 2

Comparison of algorithms for predicting peptide efficacy

We tested 5 methods for their ability to predict peptide efficacy. A complete description of each method can be found in Materials and Methods. The methods were evaluated on three datasets that measured the ability of peptides to inhibit endothelial cell proliferation (A), migration (B), and adhesion (C). For each method and dataset, we compute LOOCV test error. To determine which methods were superior to others, we conducted t-tests for all pairs of methods based on their squared test errors. Significantly low test errors indicate better performance. The table gives the p-value associated two-tailed paired t-test. At the 0.05 level, the naive featureless method had significantly higher error than all other methods. Also, the non-linear Lasso method had significantly less error than the Lasso method.

(A) Proliferation Models	Lasso	non-linear Lasso	local Lasso	local non-linear Lasso	featureless
Lasso	-	0.004	0.580	0.218	0.003
non-linear Lasso		-	0.496	0.815	0.001
local Lasso			-	0.343	0.018
local non-linear Lasso				-	0.007
featureless					-
(B) Migration Models	Lasso	non-linear Lasso	local Lasso	local non-linear Lasso	featureless
Lasso	-	0.034	0.340	0.000	0.000
non-linear Lasso		-	0.179	0.651	0.000
local Lasso			-	0.034	0.000
local non-linear Lasso				-	0.000
featureless					-
(C) Adhesion Models	Lasso	non-linear Lasso	local Lasso	local non-linear Lasso	featureless
Lasso	-	0.033	0.205	0.098	0.017
non-linear Lasso		-	0.862	0.432	0.004
local Lasso			-	0.740	0.005
local non-linear Lasso				-	0.007
featureless					-

LRRF**XNINN*	0.011	LRRF F**NINNVXN*	0.006
LRRF***XNINN N*	0.011	LRRF F**SNINNVXN*	0.006
LRRF***XNINNV N*	0.011	LRRF F**SNINNVSN*	0.006
LRRF***XNINNVXN*	0.011	LRRF T F**ININNVSN*	0.002
LRRF****XNINNVXN*	0.011	LRRFST F**ININNVSN*	0.002
LRRF****XNINNVXN**	0.011	LRRFST FM**NINNVSN*	0.002
LRRF**M**XNINNVXN*	0.009	LRRFST FMF**NINNVSN*	0.002
LRRF**MP XNINNVXN*	0.009	LRRFS**FMF**NINNVSN*	0.001
LRRF**MP**XNINNVXN*	0.003	LRRFS**FMF**NINNVSN*	0.001
LRRF**MP***XNINNVXN*	0.003	LRRFS***FMF**NINNVSN*	0.001
LRRF**MP****XNINNVXN*	0.003	LRRFS***FMF**NINNVXN*	0.001
LRRF**MP*****XNINNVXN*	0.002	LRRF****FMF**NINNVXN*	0.001
LRRF**MP*****XNINNVXN**	0.002	LRRFS***AMF**NINNVXN*	0.001
LRRF**MP*****XNINNVXN**	0.001	LRRFST**AMF**NINNVXN*	0.001
LRRF**MP*****NINNVXN*	0.001	LRRFST**AMF**NINNVXN*	0.001
LRRFSTAPPEMFXXNINNVXNF	weights	LRRFSTAPPEMFXXNINNVXNF	weights