

Computacion Avanzada en la Prediccion de la Estructura de Proteinas

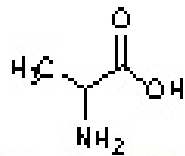
Natalio Krasnogor
www.cs.nott.ac.uk/~nxk

Proteínas a Vuelo de Pajaro

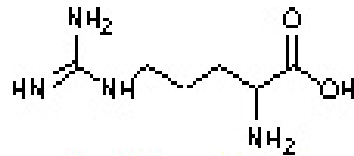
- Las proteínas son **combinaciones** de los 20 aminoácidos. A partir de estos objetos! **formadas** por (algunos) de los 20 aminoácidos.
- Se conectan formando una **cadena lineal**. Queremos predecir estos objetos
- Esta **cadena lineal** o secuencia se llama estructura primaria (compuesta por un string en el alfabeto de tamaño 20).
- La estructura primaria forma estructuras locales secundarias
- Las estructuras secundarias dan lugar a la **estructura terciaria**.

Schematic diagrams of the 20 amino acids

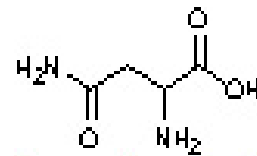
(picture taken from www.chemistry.pomona.edu)



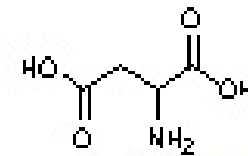
Alanine (Ala)



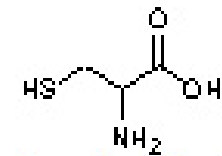
Arginine (Arg)



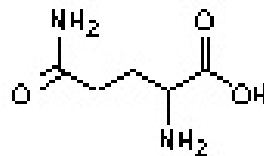
Asparagine (Asn)



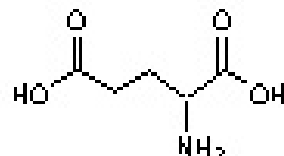
Aspartic Acid (Asp)



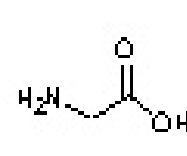
Cysteine (Cys)



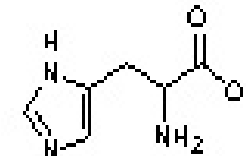
Glutamine (Gln)



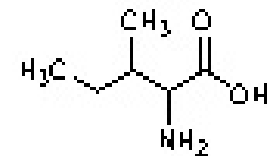
Glutamic Acid (Glu)



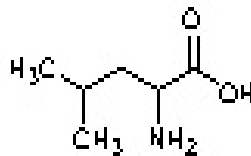
Glycine (Gly)



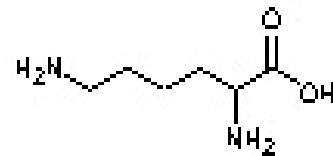
Histidine (His)



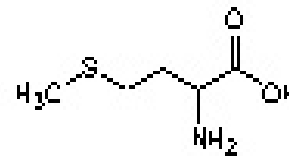
Isoleucine (Ile)



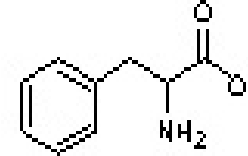
Leucine (Leu)



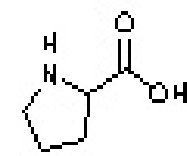
Lysine (Lys)



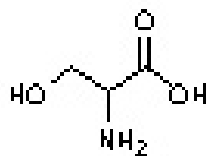
Methionine (Met)



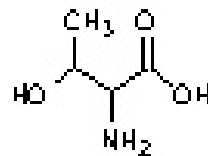
Phenylalanine (Phe)



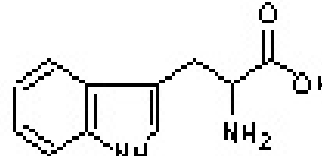
Proline (Pro)



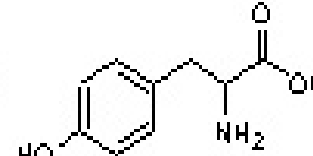
Serine (Ser)



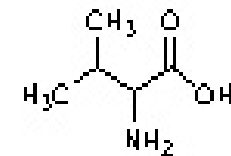
Threonine (Thr)



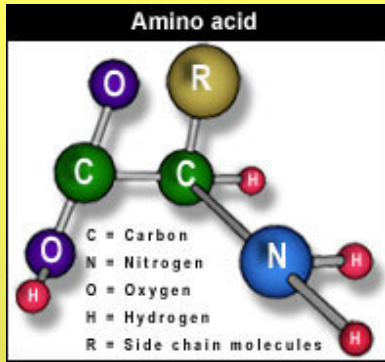
Tryptophan (Trp)



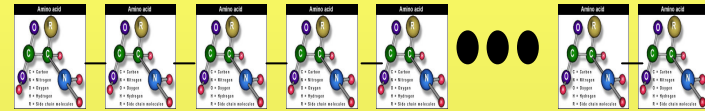
Tyrosine (Tyr)



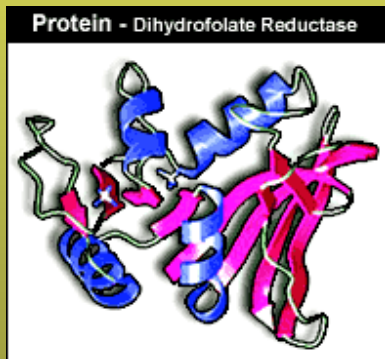
Valine (Val)



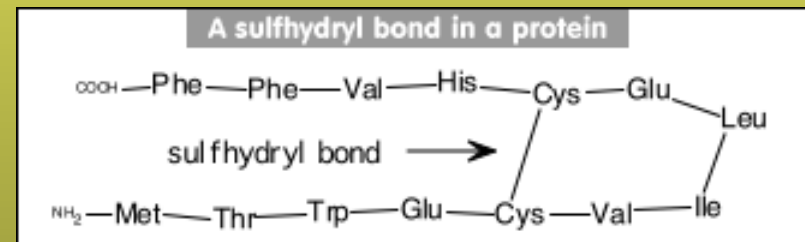
Estructura primaria



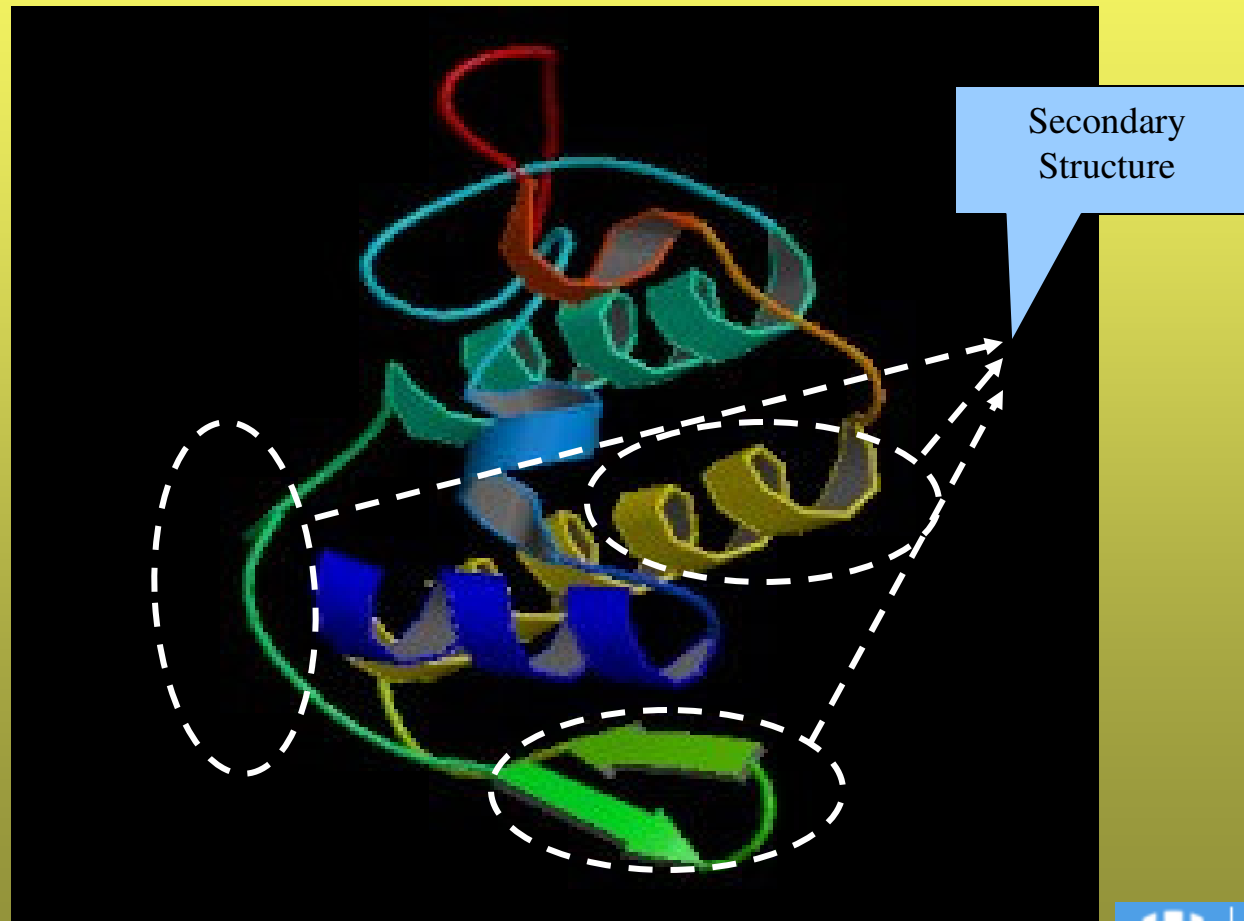
Estructura Secundaria



Estructura terciaria



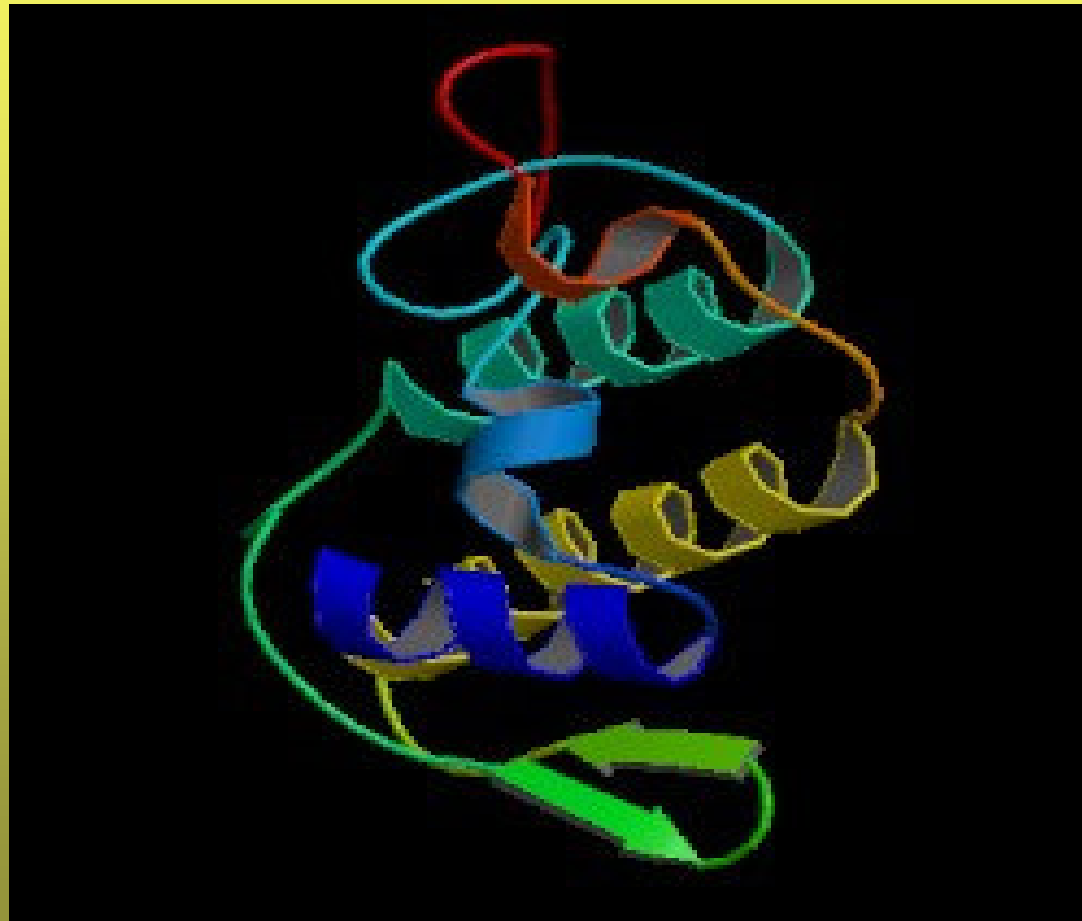
Example: 1CLP



Example: 1CLP

Tertiary
Structure

Native
State



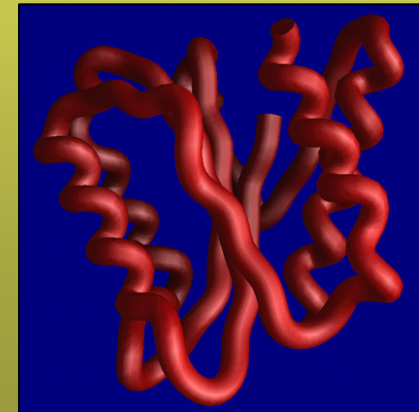
El role de las proteínas en la vida

- **Proteínas estructurales:** los bloques constructivos del organismo, ej. colageno, unias, pelo, etc.
- **Enzimas:** maquinaria biological que regulan una multitud de funciones metabolicas. Generalmente las enzimas son muy especificas y catalizan un solo tipo de reaccion, pero pueden jugar algun papel en mas de una reaccion.
- **Proteínas de las membranas:** son las “amas de casa” de las celulas, ej. Mediante la regulacion del volumen celular, extraccion y concentracion de pequenas moleculas del ambiente extracelular y la generacion gradientes ionicos essentielles para el funcionamiento muscular y nervioso (la bomba de sodio/potasio es un ejemplo)

Problema de Prediccion de la Estructura Terciaria de Proteinas:

- INPUT: una secuencia de aminoacidos
- OUTPUT: la estructura terciaria que esta adopta

```
MKIVYWSGTGNTK  
MAELIAKGIIESGK  
DVNTINVSDVNI  
DELLNEDILILGCS  
AMGDEVLEESEFEP  
FIEEISTKISGK  
KVALFGSYGWGDGK  
WMRDFEERMNGYGC  
VVVETPLIVQNE  
PDEAEQDCIEFGKK  
IANI
```

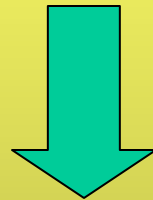


- No existe consenso entre los científicos de cual es el modelo atómico adecuado para predecir el plegado de una proteína
- Aun cuando se conociese ese modelo, hacer una simulación átomo-por-átomo del plegado sería muy costoso computacionalmente
- Modelos simplificados de proteínas han impactado en nuestro entendimiento:
 - del comportamiento de las mismas
 - de los algoritmos necesarios para predecir el plegado.

Estos modelos minimalistas se han usado entre otras cosas para:

- Estudiar la naturaleza del energy landscape (NakSasSas2001)
- La unicidad del estado nativo y las secuencias degeneradas (SunBreChaDi196)
- La transición de dos estados de estructuras de tipo globins (YueDi194)
- La existencia de plegado cooperativo, la relación entre la estructura y la función (Hir99)

- Predecir la estructura de proteínas reales por medio de una combinación de información experimental sobre estructura secundaria y terciaria con modelos optimizados obtenidos de simulaciones en latices, ej. (SamXiaHuaLev99, XiaHuaLevSam2000, FeiBro2000, KihLuKolSko2001b, KolBetKihRotSko2001)



Los modelos simplificados continúan jugando un rol importante en:

- la mejora de nuestro conocimiento de los fundamentos físicos de proteínas reales
- facilitando el desarrollo de mejores algoritmos para la predicción de estructura.

El problema de predecir el estado nativo de una proteína se descompone en:

- **Modelado:** especificacion y seleccion de la funcion de energia, geometria, uso de “side-chains”, etc.
- **Optimisation:** se concentra en las tecnicas algoritmicas que son necesarias para buscar la estructura optima de una secuencia especifica asumiendo un modelo particular.

En el resto de la charla asumiremos una familia de modelos minimalistas y nos concentraremos en los algoritmos necesarios para resolverlos.

Los Modelos Utilizados

- Modelo HP (Dil85)
- Modelo Functional Model Proteins, FMP, (Hir99, BlaHir01, BlaHir03)
- Embebidos en una geometria 2d o 3d, reticular, ej., cubicos, diamante, centrado en las caras, etc.
- HP y FMP abstraen las interacciones hidrofobicas y reducen el alfabeto de 20 aminoacidos a 2: no-polares o hidrofobicos (H) y polares (P) o hidrofílicos
- Una proteina de n aminoacidos se representa por una secuencia $s \in \{H,P\}^n$ que se mapea al reticulado donde cada residuo ocupa una celda diferente. El mapeo debe ser self-avoiding.

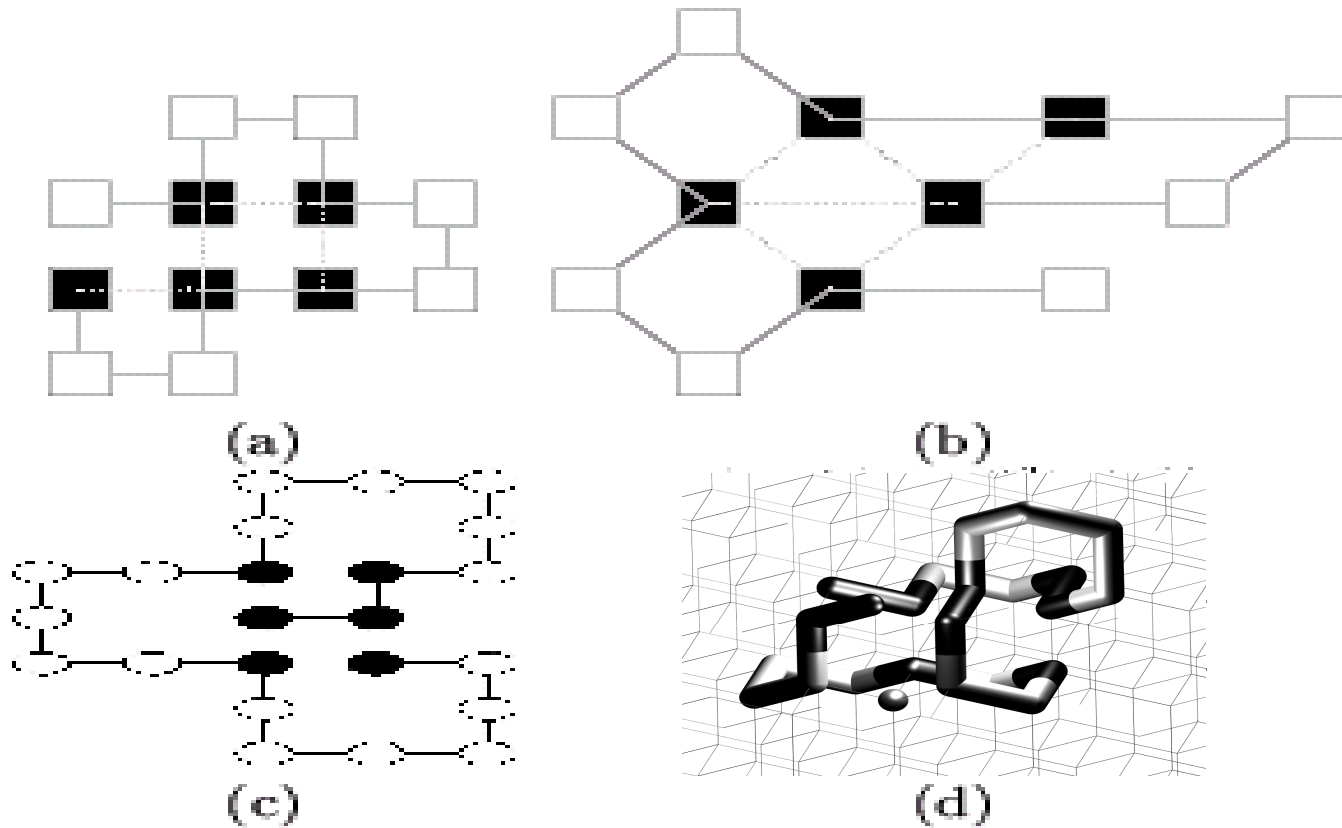


FIG. 1: *HP* protein embedded in the square lattice (a) and triangular lattice (b). Functional Model protein embedded in the square lattice (c) and diamond (3D) lattice (d). In (c) and (d) native structures are not maximally compact as they must have a “binding pocket”.

La Funcion de Energia (funcion objetivo)

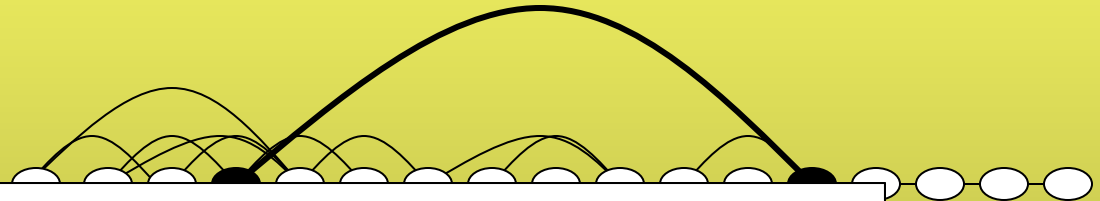
Una proteina:



Su estructura:



El mapa de contacto



$$E(s) = \sum_{i < j ; 1 \leq i, j \leq n} (\Delta_{i,j} * \epsilon_{i,j})$$

where

$$\Delta_{i,j} = \begin{cases} 1 & \text{if } i, j \text{ are in contact and } |i - j| > 1 \\ 0 & \text{otherwise} \end{cases}$$

- Una proteína del FMP debe plegarse a un estado nativo unico, tener un “binding pocket” y un “energy gap”
- Estas restricciones adicionales al model HP lo hacen un modelo mas dificil de resolver.
- Es mas probable quedarse atrapado en un optimo local pobre
- No se puede asumir compacticidad de la estructura
- Existe solo un optimo global
- Varias versiones de los modelos HP son NP-Hard (BerLei98, Cre98).

En Resumen

- 20 tipos de residuos
 - No-polares o hydrophobicos (H)
 - Polares (P) o hydrophilicocs
- Una proteina de n residuos se representa por una sequencia s
- La sequencia se mapea al latice
- Cada residuo en s ocupa una celda del latice diferente
- El mapeo debe ser “self-avoiding”

$$E(s) = \sum_{i < j ; 1 \leq i, j \leq n} (\Delta_{i,j} \epsilon_{i,j})$$

- El potencial de energia reflaja la propensidad de los H a formar close contacts
- Standard HP model
 - HP and PP assigned energy 0
 - HH contact assigned energy -1

$$\Delta_{i,j} = \begin{cases} 1 & \text{if } i, j \text{ are in contact and } |i - j| > 1 \\ 0 & \text{otherwise} \end{cases}$$

Metaheurísticas para la Predicción de Estructuras

(Glo86) define metaheurísticas como:

...a master strategy that guides and modifies other heuristics to produce solutions beyond those that are normally generated in a quest for local optimality. The heuristics guided by such a meta-strategy may be high-level procedures or may embody nothing more than a description of available moves for transforming one solution into another, together with an associated evaluation rule.

Algunas de las metaheurísticas más exitosas para estos modelos incluyen:

- Automatas Celulares (KraMarPelRiz98)
- Algoritmos Genéticos (Krasnogor99)
- Híbridos Evolutivos-Monte Carlo (LiaWon2001)
- Búsqueda por Vecindarios Variables Fuzzy (KraPel2002)
- Algoritmos Meméticos (KraBlaBurHir2002)
- Optimización por Colonia de Hormigas (ShmAguHoo02)
- Híbridos Genéticos-Búsqueda Tabu (JiaCuiShiMa03)

Principios de Diseño para Algoritmos Exitosos

Los algoritmos mas exitosos usan al menos 3 de los siguientes:

1. The use of global search mechanisms in contrast with purely local search.
2. The explicit use of diversity preserving mechanisms.
3. The use of a number of distinct move operators rather than a single move operator.
4. The explicit use of mechanisms to jump-out of poor local optima and/or traverse large neutral plateaus.
5. The use a “memory” of visited configurations.

El primero principio de disenio es generalmente implementado por medio de un (multi)set de soluciones.

Teniendo varias soluciones es mas dificil quedarse atrapado en optimos locales pobres.

Sin embargo, este principio se relaciona con el segundo, la preservacion de la diversidad. En el peor caso todo el (multi)set representaria la misma solucion lo cual implicaria una perdida de recursos computacionales.

Las soluciones exploradas deben ser diversas de tal manera que representen varias regiones distintas del espacio de busqueda.

En relacion a la diversidad y a los optimos locales se introduce el tercer principio de disenio: el uso de varios operadores de busqueda.

Si el algoritmo tiene a su disposición una variedad de “move operators” entonces es más difícil que pierda diversidad o que se quede atrapado en óptimos locales pobres.

En general los óptimos locales de varios tipos de operadores son valientes

Sin embargo, este mecanismo por sí solo no puede garantizar no quedarse atrapado en óptimos locales.

El cuarto principio de diseño apunta a ello. Hay que incluir mecanismos que detecten óptimos locales (o mesetas) para poder escaparse de ellos.

Finalmente, el quinto principio apunta a mantener alguna forma de memoria de configuraciones visitadas.

- principios no ortogonales
- pueden ser conflictivos

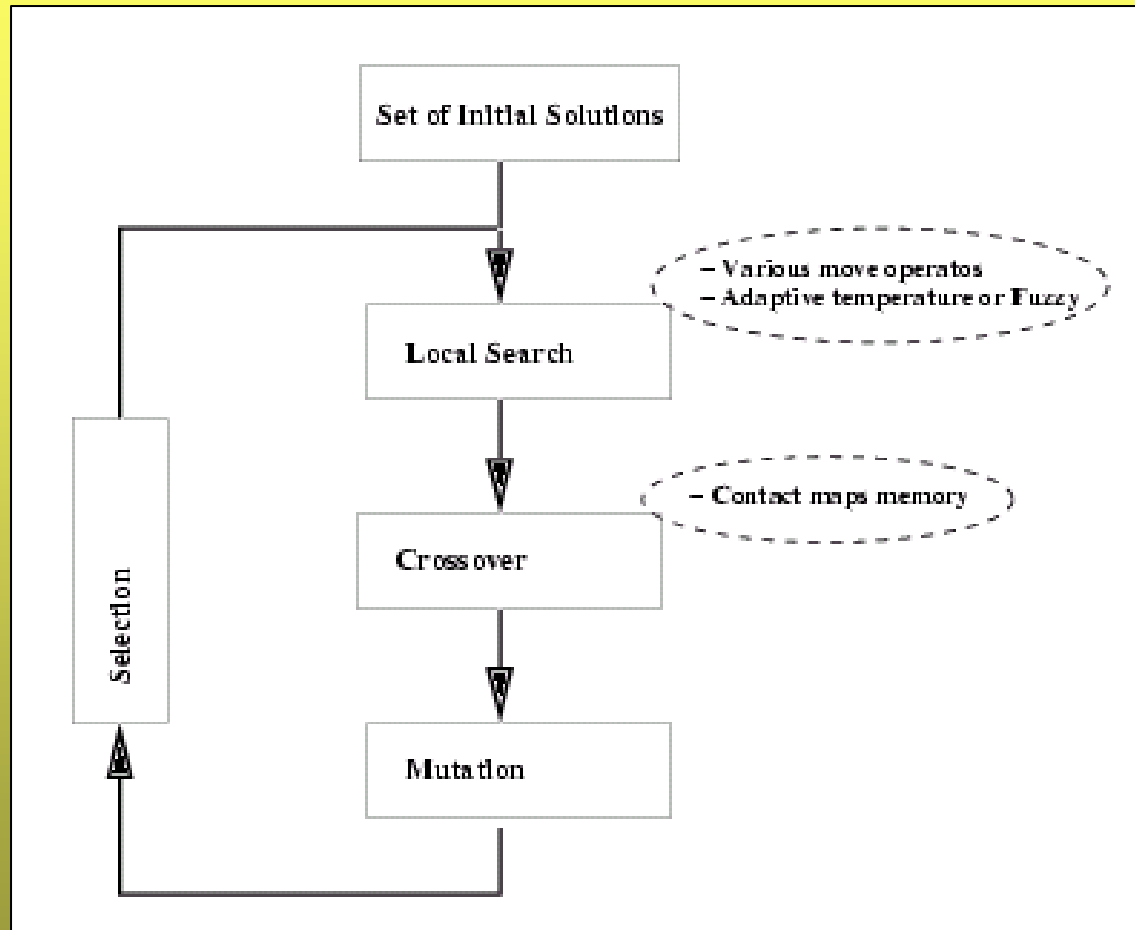
Ejemplos de Algoritmos Recientes que Emplean Algunos de Estos Principios

Hibridos Evolutivos-Monte Carlo

(KraSmi99b) utiliza 1 (poblacion) , 2 (diversificacion) & 4 (evitar optimos locales)

Usa un hibrido (algoritmo memetico) Evolutivo-Monte Carlo

El Evolutivo implementa el principio 1 y el MC modificado el 2&4



El MC modificado sirve de buscador local cuando la poblacion es diversa y de diversificador cuando la poblacion del Memetico es homogenea

La temperatura del MC se define antes de cada ejecucion como:

$$T = \frac{1}{|E(v_h) - E(v_l)|}$$

Y se acepta una nueva solucion de acuerdo a :

$$\text{random}(0, 1) < \exp^{-k * \frac{\delta_{w, \text{new}}}{T}}$$

cuando la poblacion es diversa, $T \rightarrow 0$, acepta solo mejoras

cuando la poblacion es Homogenea, $T \rightarrow \infty$, acepta mejores/peores soluciones

(LiaWon2001) tambien usa un Memetico con MC e implementan los principios 1, 2 & 3:

- mantienen una poblacion de soluciones (1)
- usan varios operadores de busqueda que le permite al algoritmo inducir un paisaje de busqueda rico (3)
- incluyen un operador de intercambio para las temperaturas del MC lo que le permite al algoritmo mantener la diversidad (2) ya que puede funcionar como muy explorativo (altas T) a muy explotativo (bajas T):
 - parallel tempering (GeyTho95, YanPab99)
 - single chain Metropolis-Hastings algorithm.

Colonia de Hormigas

En (ShmAguHoo02) se implementa una CH

A diferencia de los Memeticos con MC, la CH mantiene una poblacion de hormigas (no soluciones explicitas) (1).

Cada hormiga tiene asociado un vector de probabilidades de instanciar una solucion particular

El principio 2 (diversidad) se implementa con un operador de re-inicio que realoca probabilidades a estos vectores

Ademas la CH utiliza dos operadores de busqueda local (principio 3)

Hibridos Evolutivos- Búsqueda de Operador Variable Fuzzy

En (KraPel2002) los autores usan búsqueda de vecindarios variables fuzzy (FANS).

Los principios 1, 3 & 4 son empleados.

Un algoritmo evolutivo implementa el primer principio y se hibridiza con FANS.

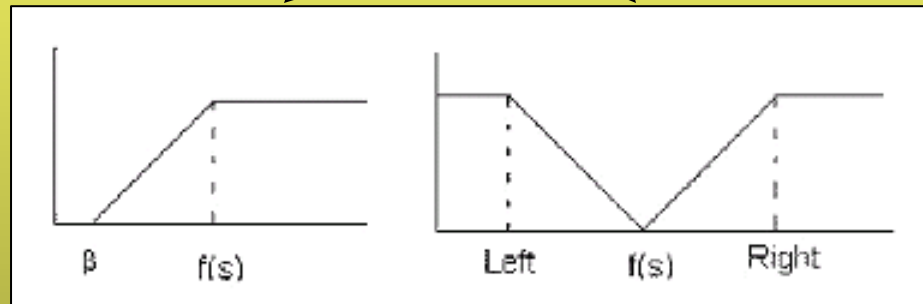
FANS se usa para enfocar la búsqueda en buenas regiones del espacio de búsqueda inducido por varios operadores (3)

Un criterio fuzzy se usa de manera semejante al MC modificado de manera que de vez en cuando se aceptan soluciones con energías más altas (4).

$$\mu(E(v), E(v_{new})) \geq \lambda$$

funcion de pertenencia al conjunto de soluciones
“acceptables”

Dos conceptos distintos de aceptabilidad



Promueve mejoras

Promueve diversidad

$$\mu_1(v, v_{new}) = \begin{cases} 0.0 & \text{if } E(v_{new}) > \beta \\ \frac{(E(v_{new}) - \beta)}{(E(v) - \beta)} & \text{if } \beta \geq E(v_{new}) \geq E(v) \\ 1.0 & \text{if } E(v_{new}) < E(v) \end{cases}$$

$$\mu_2(v, v_{new}) = \begin{cases} 1.0 & \text{if } Right \geq E(v_{new}) \geq Left \\ \frac{(E(v) - E(v_{new}))}{(E(v) - Left)} & \text{if } Left > E(v_{new}) \geq E(v) \\ \frac{(E(v_{new}) - E(v))}{(Right - E(v))} & \text{if } E(v) > E(v_{new}) \geq Right \end{cases}$$

Algunos Resultados (1)

Static Move Operators	#O/#R	Mean FHT
GA (no move operators)	0/10	-
MA with Macro Mutation (r=4)	2/10	27.5
MA with Macro Mutation (r=8)	3/10	53.3
MA with Macro Mutation (r=16)	2/10	43.0
MA with Reflect (r=4)	3/10	20.6
MA with Reflect (r=8)	1/10	79.0
MA with Reflect (r=16)	1/10	45.0
MA with Stretch (r=4)	0/10	-
MA with Stretch (r=8)	0/10	-
MA with Stretch (r=16)	0/10	-
MA with Pivot	5/10	27.0
MultiMeme (all local searchers)	8/10	16.87

Principio 3: cero, uno o varios operadores
modelo hp 2d, rectangular

Algunos Resultados (2)

#	Sequence	Size	Opt.	MMA
1	$HPHRPHHRPHRPHHRPHRPH$	20	-9	-9
2	$PPRHHRRHHRRPPRRHHHHHHHRPHHRPPRRHHRRPHRPP$	36	-14	-14
3	$H^2(PH)^4H^3RHR^3HR^3HR^4HR^3HR^3HRH^4(PH)^4H$	50	-21	-21
4	$H^{12}(PH)^2(P^2H^2)^2(PRH)^2(HR^2H)^2(P^2H^2)^2P^2(HR)^2H^{12}$	64	-42	-39
5	$HPHRPHHRPHRPHHRPHRPH$	20	-9	-9
6	$PRHRPHHRPPRRHHRRPPRRHHRRPPRRHH$	25	-8	-8
7	$(P^2H)^2HR^2H^2P^5H^{10}P^6(H^2P^2)^2HR^2H^5$	48	-22	-22(-23)
8	$RHRPHRHHHRHHRH^5$	18	-9	-9
9	$HRHRHHHRPPRRHHHRPHH$	18	-8	-8
10	$HHRRPPRRHHRRPPRRPHR$	18	-4	-4
11	$HHHRPHRPHRPHRPHRPHRPH$	20	-10	-10
12	$PRH^3RH^8P^3H^{10}RHR^3H^{12}P^4H^6RHHRPH$	60	-34	-34(-35)

Principio 4: Fuzzy logic para salir de optimos locales.
Nuevos optimos en el modelo HP, 2d

Algunos Resultados (3)

#	Sequence	Size	Opt.	MMA
1	<i>HHPRHPRHPRHP</i>	14	-11	-11
2	<i>HHPRHPRHPRHPH</i>	14	-11	-11
3	<i>HHPRHPRHPRHPRH</i>	16	-11	-11
4	<i>HHPRHPRHPRHPRHP</i>	16	-11	-11
5	<i>HHPRHPRHPRHPRHPH</i>	17	-11	-11
6	<i>HHPRHPRHPRHPRHHPH</i>	17	-17	-17
7	<i>HHPRHPRHPRHPRHPRHH</i>	20	-17	-17
8	<i>HHPRHPRHPRHPRHPRHH</i>	20	-17	-17
9	<i>HHPRHPRHPRHPRHPRHH</i>	21	-17	-17
10	<i>HHPRHPRHPRHPRHPRHH</i>	21	-17	-17
11	<i>HHPRHPRHPRHPRHPRHH</i>	21	-17	-17
12	<i>HHPRHPRHPRHPRHPRHH</i>	22	-17	-17
13	<i>HHPRHPRHPRHPRHPRHHH</i>	23	-25	-25
14	<i>HHPRHPRHPRHPRHPRHHH</i>	24	-17	-16
15	<i>HHPRHPRHPRHPRHPRHHH</i>	24	-25	-25
16	<i>HHPRHPRHPRHPRHPRHHH</i>	24	-25	-25
17	<i>HHPRHPRHPRHPRHPRHHH</i>	30	-25	-24
18	<i>HHPRHPRHPRHPRHPRHHH</i>	30	-25	-24
19	<i>HHPRHPRHPRHPRHPRHHH</i>	37	-29	-26

Principio 3: MC modificado en el reticulado triangular, modelo hp 2d.

Algunos Resultados (4)

Un algoritmo memetico que usa multiples operadores, el MC modificado y memoria en forma de contact maps (los 5 principios):

```
Reproduction Stage(P):  
Begin  
  /* The parent population is stored in P */  
   $O = \emptyset$ ;  
  /* O, the offspring set, is initialised */  
  Generate the contact map memory for  $P \rightarrow CMM$ .;  
  While ( Not enough offspring generated ) Do  
    Select two parents to mate,  $p_1, p_2$ ;  
     $crossover(p_1, p_2) \rightarrow offspring$ ;  
    If ( $compatible(offspring, CMM)$ ) Then  
       $O = O \cup offspring$ ;  
      Inherit move operator to offspring using SIM;  
    Fi  
  Od  
  Return  $O$ ;  
End.
```

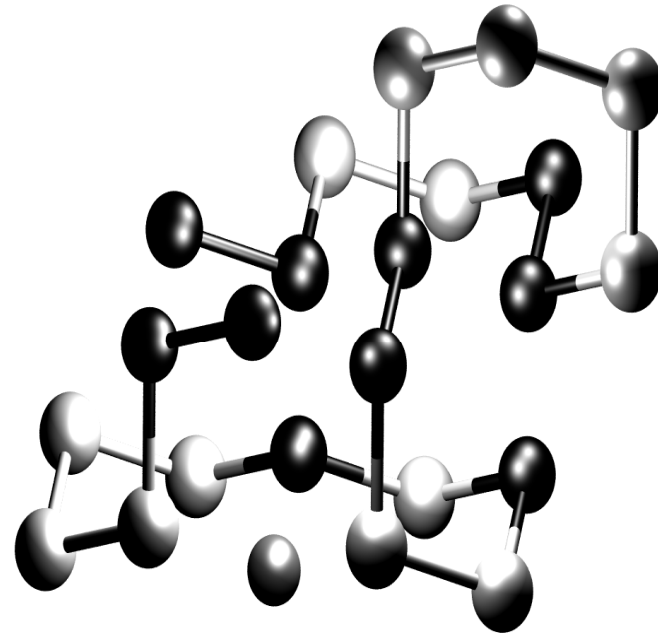
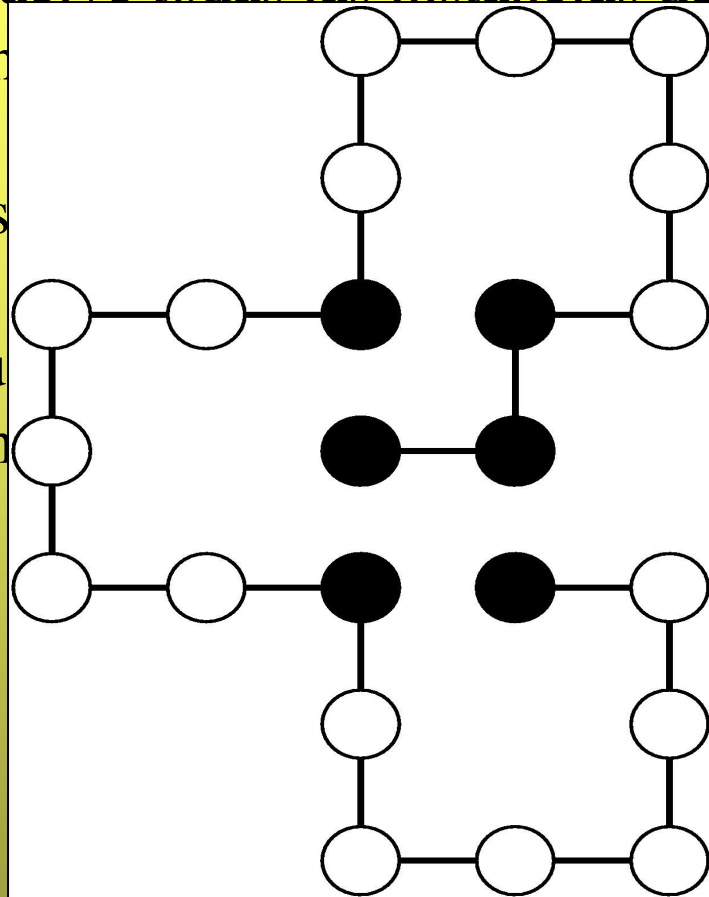
Resuelve todas las instancias del HP en reticulados rectangulares y

trian

Mas

resu

mon



el

Conclusiones

El desarrollo de algoritmos robustos en modelos simplificados son un “stepping-stone” para alcanzar algoritmos efectivos en modelos mas realistas que no pueden ser resueltos por “threading” o “homology”

Varios de los algoritmos mas recientes presentados en CASP utilizan un sampleo de estructuras en latices

En esta charla destile varios principios de disenio importantes para algoritmos en modelos simplificados. Es improbable que un algoritmo que no utilice esto principios pueda mejorar los resultados actuales

Nuestro algoritmo memetico (de 5 principios) es el mas robusto ya que resuelve instancias de los modelos HP y FMP en dos y tres dimensiones y en varias geometrias (rectangular, triangular, fcc)

Veremos luego que estos modelos,
extremadamente simple,
no nos alejan demasiado de los problemas
reales!