































Model for a ribosome binding site (based on ~300 known RBS's)

	1	2	3	4	5
Т	0.161	0.050	0.012	0.071	0.115
С	0.077	0.037	0.012	0.025	0.046
A	0.681	0.105	0.015	0.861	0.164
G	0.077	0.808	0.960	0.043	0.659

Nucleic Acids Research, 1998, Vol. 26, No. 4

1107-1115

How well does it do on well-characterized genomes?						
Genome	Genes annotated	Genes predicted	Exact prediction (%)	Missing genes (%)	Wrong genes (%)	
A.fulgidus	2407	2530	73.1	10.8 (2.0)	15.1	
B.subtilis	4101	4384	77.5	3.6 (2.8)	9.8	
E.coli	4288	4440	75.4	5.0 (2.7)	8.2	
H.influenzae	1718	1840	86.7	3.8 (3.2)	10.2	
H.pylori	1566	1612	79.7	6.0 (4.4)	8.7	
M.genitalium	467	509	78.4	9.9 (1.7)	17.3	
M.jannaschii	1680	1841	72.7	4.6 (0.8)	12.9	
M.pneumoniae	678	734	70.1	7.8 (4.1)	13.6	
M.thermoauthotrophicum	1869	1944	70.9	5.0 (3.5)	8.6	
Synechocystis	3169	3360	89.6	4.0 (1.5)	9.4	
Averaged	21 943	23 194	78.1	5.4 (2.7)	10.4	
But this was a long time ago!						



















GENSCAN, when it was first developed						
		Accuracy		Accuracy		
		per base		per e	per exon	
Program	Sequences	Sn	Sp	Sn	Sp	
GENSCAN	570 (8)	0.93	0.93	0.78	0.81	
FGENEH	569 (22)	0.77	0.88	0.61	0.64	
GeneID	570 (2)	0.63	0.81	0.44	0.46	
Genie	570 (0)	0.76	0.77	0.55	0.48	
GenLang	570 (30)	0.72	0.79	0.51	0.52	
GeneParser2	562 (0)	0.66	0.79	0.35	0.40	
GRAIL2	570 (23)	0.72	0.87	0.36	0.43	
SORFIND	561 (0)	0.71	0.85	0.42	0.47	
Xpound	570 (28)	0.61	0.87	0.15	0.18	
GeneID+	478 (1)	0.91	0.91	0.73	0.70	
GeneParser3	478 (1)	0.86	0.91	0.56	0.58	
				J. Mol. Bio	ol. (1997) 268 , 78–94	





How well do we know the genes now?In the year 2000"Over 95% of the coding nucleotides ... were correctly
identified by the majority of the gene finders.""...the correct intron/exon structures were predicted for >40%
of the genes."Most promoters were missed; many were wrong."Integrating gene finding and cDNA/EST alignments with
promoter predictions decreases the number of false-positive
classifications but discovers less than one-third of the
promoters in the region."

Genome Research 10:483-501 (2000

How well do we kno	now?			In the year 2006		
	Table 3					
	Summary of programs used t	o determine predictions submitt	ed for each EGASP category			
EGASP: the	Submission category	Program	Affiliation	Reference	Assessment	
Project	I (AUGUSTUS-any) 2 (AUGUSTUS-abinit) 3 (AUGUSTUS-EST) 4 (AUGUSTUS-dual)	AUGUSTUS	Georg-August-Universität, Göttingen	[58]		
	T	FGENESH++	Softberry Inc.	[56]		
	1	JIGSAW	The Institute for Genomic Research (TIGR)	[59]		
= scientists i	I (PAIRAGON-any) 3 (PAIRAGON+NSCAN_EST)	PAIRAGON and NSCAN_EST	Washington University, Saint Louis (WUSTL)	[57]	SP)to	
prodict gone	2	GENEMARK.hmm	Georgia Institute of Technology	[60]	ara tham to	
I predict gene	2	GENEZILLA	TIGR	[81]	are them to	
	3	ACEVIEW	National Center for Biotechnology Information (NCBI)	[52]		
experimenta	3	ENSEMBL	The Wellcome Trust Sanger Institute (WTSI) and European Bioinformatics Institute (EBI)	[64]		
	• vve	EXOGEAN	Ecole Normale Superieure, Paris	[62]		
	discussed	EXONHUNTER	University of Waterloo	[63]		
10 groups	 aiscussea 	ACESCAN*	Salk Institute	[82]		
TA groups	those	DOGFISH-C	WTSI	[67]		
	these	NSCAN	WUSTL	[57]		
26 programs	apriliar	SAGA	University of California at Berkeley	[66]		
30 programs	earner	CENED UID	WOSTL - EDI	[63]		
1 0	2	SCP2 LU2	Médica Barcelona	-		
	6	ASPICT	Liniversità degli Studi di Milano	(83)		
	6 (AUGUSTUS-exon)	AUGUSTUS	Geore-August-Universität, Göttingen	[58]		
	6	CSTMINER [‡]	Università degli Studi di Milano	[84]		
	6	DOGFISH-C-E [§]	WTSI	[67]		
	6	SPIDA	EBI	[85]		
	6	UNCOVER§	Duke University	[86]		
	L	CCDSGene	UCSC tracks [7]	[55]		
	I.	KNOWNGene		[54]		
	1	REFSEQ (REFGene)		[4]		
	2	GENEID		[19]		
	2	GENSCAN		[18]		
	3	ACEMBLY		[52]		
	2	EDISEMPI (EDISCome)		[53]		
	3	MGCGene		[5]		
	4	SGP2		[9]		
	4	TWINSCAN		[12,13]		
		CODING 20050607	GENCODE annotation	[33]		
	-	GENES 20050607			Genome Biology 2006, 7(Suppl 1):S	























What about the current state of prokaryote gene models?

- "We applied AssessORF to compare gene predictions offered by GenBank, GeneMarkS-2, Glimmer and Prodigal on genomes spanning the prokaryotic tree of life.
- Gene predictions were 88–95% in agreement with the available evidence, with Glimmer performing the worst but no clear winner.
- All programs were biased towards selecting start codons that were upstream of the actual start."

Bioinformatics, 36(4), 2020, 1022–1029





