

Gcm1, a mammalian homolog of *Drosophila* Glial Cells Missing

Yelena Altshuller^a, Neal G. Copeland^b, Debra J. Gilbert^b, Nancy A. Jenkins^b,
Michael A. Frohman^{a,*}

^aDepartment of Pharmacological Sciences and Institute for Cell and Molecular Biology, State University of New York,
Stony Brook, NY 11794-8651, USA

^bMammalian Genetics Laboratory, ABL-Basic Research Program, NCI-Frederick Cancer Research and Development Center,
Frederick, MD 21702, USA

Received 12 June 1996

Abstract Differentiation of glia (astrocytes and oligodendrocytes) in *Drosophila* requires the gene *glial cells missing* (*gcm*), which controls lineage determination. In the absence of *gcm*, neuroglia progenitors exclusively differentiate into neurons, instead of into both neurons and glia. In contrast, ectopic overexpression of *gcm* causes uniform differentiation of the neuroglia progenitors into glia. Glial and neuronal cells in vertebrates similarly derive from neuroblast progenitors. To investigate vertebrate glial formation, we have identified, cloned, and chromosomally mapped a mammalian *gcm* homolog. Mouse *Gcm1* demonstrates extensive similarity to *Drosophila gcm* but is expressed at very low levels during neuro- and gliogenesis.

Key words: *Gcm*; Glia; Glial cells missing; Gliogenesis; mGCM1; Neurogenesis

1. Introduction

In the adult mammalian brain, astrocytes and oligodendrocytes, known collectively as glia, constitute 90% of the cells present (reviewed in [1,2]). Glia arise during development from the same lineage as neurons. Embryonic ectoderm differentiates in part into neuroectoderm during gastrulation, followed by formation in succession of the neural plate, neural folds, the neural tube, and then specific regions along the anteroposterior axis including the spinal cord, and different regions of the brain. During each of these steps, the cells that populate the neural field constitute multi-potential neuro-glia progenitors, commonly called neuroblasts. Glia first start to differentiate in mid-embryogenesis with the emergence of 'radial glia' which act to form a scaffold over which differentiating neurons subsequently migrate from the center of the neural tube to reach their final peripheral destinations. A much larger wave of gliogenesis occurs near birth and during the first 2 post-natal weeks, as the nervous system matures. Lineage marker analysis indicates that multipotential cells capable of giving rise to both neurons and glia exist at least until birth, although they are easier to demonstrate at earlier stages of development. Nothing is known about the molecular signals that control the choice of a progenitor cell to differentiate into one lineage or the other.

glial cells missing (*gcm*) was identified as a gene that regulates this decision in *Drosophila* [3,4]. In the absence of *gcm*,

neuroglia progenitors exclusively differentiate into neurons, instead of into both neurons and glia. In contrast, ectopic overexpression of *gcm* causes uniform differentiation of the neuroglia progenitors into glia (reviewed in [5]). *gcm* encodes a novel protein with limited homology to any other known proteins. *gcm* encodes a nuclear localization sequence, localizes to the nucleus, and causes changes in transcription, suggesting that it might function either as a novel transcription factor or as a transcription accessory factor. The protein also contains three copies of a four amino acid repeat (MPVP), a PEST sequence suggesting that it is rapidly degraded, a highly basic cysteine-rich region, and homopolymeric stretches of glutamine, tyrosine, serine, and glycine. Finally, the mRNA itself encodes an RNA instability motif, AUUUA, in its 3' untranslated region. *gcm* is expressed in *Drosophila* only transiently, during specific early steps of gliogenesis.

The similar origin of glia from neuroglial progenitors in *Drosophila* and vertebrates [1] prompted us to search for a mammalian homolog for *gcm*.

2. Materials and methods

2.1. cDNA isolation and characterization

RNA was prepared and reverse transcribed using a dT-tailed adapter as previously described [6]. Amplification using degenerate primers and subsequent RACE cDNA cloning was carried out as previously described [6,7]. cDNA clones were sequenced on both strands using Sequenase kits (US Biochemical).

2.2. Expression analysis

Analysis of *Gcm1* RNA levels using PCR and RNA in situ hybridization were performed as previously described using [³²P]UTP-labeled probes [7].

2.3. Interspecific mouse backcross mapping

Interspecific backcross progeny were generated by mating (C57BL/6J × *M. spretus*)F₁ females and C57BL/6J males as described [8]. A total of 205 N₂ mice were used to map the *Gcm1* locus (see text for details). DNA isolation, restriction enzyme digestion, agarose gel electrophoresis, Southern blot transfer and hybridization were performed essentially as described [9]. Both coding region and UTR probes were used. The data generated with a ~200 bp 5' UTR probe are shown in Fig. 4. With the 5' UTR probe, fragments of 13.0 and 5.3 kb were detected in *Bgl*I digested C57BL/6J DNA and fragments of 19.5 and 15.0 kb were detected in *Bgl*I digested *M. spretus* DNA. The presence or absence of the *M. spretus*-specific fragments, which segregated independently, was followed in backcross mice.

A description of the probes and RFLPs for the loci linked to the *Gcm* loci has been reported previously [10]. Recombination distances were calculated as described [11] using the computer program SPRETUS MADNESS. Gene order was determined by minimizing the number of recombination events required to explain the allele distribution patterns.

*Corresponding author. Fax: (1) (516) 444-3218.
E-mail: MichaelA@pharm.som.sunysb.edu

ACAGAACTTGCAGGACCTGAGCCAGAGCTGACACTTGTATGACACCTGGACAAGAAGCTGACAAGAAGGACAGACTCGGGAGCAGTGAAG	90
CGCGACAGGCTTTGAAAAACAAGCCCTTCAGCTGCTCTTGTGGCCCGAGTTCTGAGACTCCAAATAGCTCTGACAACTGGAGTAGAGAAG	180
<u>AGCCTGTGTTGAGCAGACCTTATCATGGAACGGACGACTTTGATCCTGAAGACAAAGAGATACTGAGCTGGGACATTAACGATGTGAAA</u>	270
M E L D D F D P E D K E I L S W D I N D V K	22
CTGCCTCAGAACGTGAAAACGACTGACTGGTTCCAGGAGTGGCCGGACTCCTACGTGAAACACATCTACAGCTCGGACGACAGGAACGCA	360
L P Q N V K T T D W F Q E W P D S Y V K H I Y S S D D R N A	52
CAGCGCCACCTGAGCAGCTGGGCCATGCGCAACACCAACAACCACAACCTCCCGGATCCTCAAGAAGTCATGCCTGGGGGTAGTGGTGTG	450
Q R H L S S W A M R N T N N H N S R I L K K S C L G V V V C	82
AGCAGGGACTGCTCCACAGAGGAAGCCGCAAGATTTCCTGAGACCCGCCATCTGTGACAAAGCCAGACAGAAGCAGCAGAGGAAAAGC	540
S R D C S T E E G R K I Y L R P A I C D K A R Q K Q Q R K S	112
TGTCCCAACTGCAATGGACCCCTGAAGCTTATTCCTGCCGAGGCCACGGCGGCTTCCCGGTACCAACTTCTGGAGGCACGACGGACGC	630
C P N C N G P L K L I P C R G H G G F P V T N F W R H D G R	142
<u>TTTCATCTTTTCCAGTCCAAAGCGAGCATGATCATCCAGGCCGGAGACCAAGCTGGAAGCAGAGGCAAGAAGGCCATGAAGAAAGTG</u>	720
F I F F Q S K G E H D H P R P E T K L E A E A R R A M K K V	172
CACATGGCCTCTGCCTCCAACTCCTTACGGATGAAGGGGAGGCCAGCAGCGAAGGCGCTTCTGCTGAAATCCCGAGTCAGGGAAGTTTA	810
H M A S A S N S L R M K G R P A A K A L P A E I P S Q G S L	202
CCTTTAACTTGGTCATTCAGGAAGCGTCCAACCTGCCCGGCCTTACAGCACACCTTTAATAGCTAACGCCCCCAGCAGAATCCCTG	900
P L T W S F Q E G V Q L P G T Y S T P L I A N A P Q Q N S L	232
AATGATTGCTTATCCTTCCCCAAGAGTTATGATTTGGGGGGAAGCACTGAGCTGGAAGATCCAACTTCCACCTTAGATTCCATGAAGTTC	990
N D C L S F P K S Y D L G G S T E L E D P T S T L D S M K F	262
TATGAGAGATGCAAACTTCTCCAGCAGTAGGATCTACGGCAGTGAAGAGCAGTTTCAGCCTCCTGTCCCTGGGACGTATGGAGACTACGAA	1080
Y E R C K F S S S R I Y G S E E Q F Q P P V P G T Y G D Y E	292
GACCTGCAAACTTGAATAAAAAATGTGCCTTAGGGAGAAATCCCTCCGATGACATCTACTATCCAGCCTATCCTCTGCCTGTGGCCAGC	1170
D L Q T W N K N V A L G R N P S D D I Y Y P A Y P L P V A S	322
TGGCCCTACGACTACTTTCCCTCCAGAACTCTTTGGAGCACTTACCCAGCAAGTTCCATCAGAACCCCTGCAGCTCAACCAGGCTGT	1260
W P Y D Y F P S Q N S L E H L P Q Q V P S E P P A A Q P G C	352
CATCCCTTGTGGTCCAATCCGGGAGGTGAGCCTTACGAAGAGAAAGTATCTGTGGATCTGAGCAGCTATGTGCCCTCCCTCACGTACCAC	1350
H P L W S N P G G E P Y E E K V S V D L S S Y V P S L T Y H	382
CCACCTCAGCAGGATCCCTTCTGCTCACCTATGGCTCTCCTACTCAGCAGCAACATGCACTGCCCGGCAAGAGCAACAGGTGGGATTTT	1440
P P Q Q D P F L L T Y G S P T Q Q Q H A L P G K S N R W D F	412
GACGAAGAGATGGCATGCATGGGCTTAGATCACTTCAACAATGAGATGCTCTAAACTTCTGTCTTTAAGATAACTCAAAGCTCCCTTT	1530
D E E M A C M G L D H F N N E M L L N F C S L R *	436
CTCTCTCCAGTCGTGATTTGTTAAAGGCTTGGCAGAATTTTCTTAGAAGACGGGTTTCAAAGCATAGCCTGAGAGGAATGTTGGAATC	1620
AATTGTGACAAGCAGTAGACAAGGTCTTCTCACTCAAGCTGTGACAATGCCAGTGTGTTTACCCTCCATAAGGAATAGAAAGTCACCGT	1710
GCCAGGTGTAATGTACAAGTCTTAATCCCGCAGCTCAGGAGACAGAGGCAGTTGGATATTTGTGAGTTCAAGGCCAGCCGAGTTTACA	1800
TACTGAGTTCCAGGAGAGCCAGGACTATGTAGAGAGACCTGTCTCAACACCTCCCCCACAAGAAGAAGAAG	1875

Fig. 1. The nucleotide sequence of *Gcm1* and conceptual translation of the open reading frame conserved with *Drosophila gcm* are shown. The short open reading frame in the 5' region is underlined. The region exhibiting conservation with *gcm* is denoted with a double underline.

3. Results and discussion

3.1. cDNA sequence of *Gcm1*

A blast search with the published *Drosophila gcm* sequence identified a human expressed sequence tag that encoded a protein with a limited length of significant homology (see Fig. 3). Degenerate primers corresponding to sequences identically conserved in both the human and fly genes were used to amplify a corresponding mouse cDNA fragment. Placenta was used as the source of RNA, since the human expressed sequence tag had been generated from a human placental library. A fragment of appropriate size was identified, cloned, and sequenced. The sequence closely matched the human sequence and was used to design RACE primers which were then employed to amplify the 5' and 3' cDNA ends. The composite sequence, denoted *Gcm1*, and the conserved trans-

lated product are shown in Fig. 1A. Amplification of the entire coding region using 5' and 3'UTR specific primers was performed to ensure that assemblage of all of the pieces generated during cloning created a genuine contiguous cellular transcript (data not shown).

The *Gcm1* transcript generated is 1860 nt in length and encodes a 436 amino acid protein. Two noteworthy features are found in the 3' UTR. First, part of the 3'UTR (nts 1685–1874) contains a moderately repetitive sequence which when used in blast searches identifies scores of related sequences in otherwise unrelated genes and genomic fragments. Confirmation that this sequence is moderately repetitive came from attempting to use this region for RNA in situ hybridization and genomic southern analyses. In both cases, unacceptably high levels of background hybridization were observed (data not shown).

```

MELDDDFPED KEILSWDIND VKLPQNVKTT DWFQEWPD SYVKHIYSSDDR 50
NAQRHLSSWA MRNTNNHNSR ILKKSC LGVV VCSRDCTEE GRKIYLRPAI 100
CDKARQKQQR KSCPCNGNPL KLIPCRGHGG FVPTNFWRHD GRFIFQSKG 150
EHDHPRPETK LEAEARRAMK KVHMASASNS LRMKGRPAK ALPAEIPSG 200
SLPLTWSFQE GVQLPGTYST PLIANAPQON SLNDCLSF PKSYDLGGSTEL 250
EDPTSTLDSM KFYERCKFSS SRIYGSEEQF QPPVPGTYGD YEDLQTWNKN 300
VALGRNPDD IYYPAYPLPV ASWPYDYFPS QNSLEHLPQQ VPSEPPAAQP 350
GCHPLWSNPG GEFYEEKVSV DLSYVPSLT YHPPQDDFL LTYGSPTQQQ 400
HALPGKSNRW DFDEEMACMG LDHFNEMLL NFCSLR* 436

```

Fig. 2. The *Gcm1* protein sequence. The region conserved with *gcm* is boxed. Conserved cysteines are denoted by asterisks. A potential PEST sequence and a single glutamine repeat are underlined.

The second unusual feature is found at the extreme 3' end of the many of the RACE clones generated, none of which encode a canonical polyadenylation signal sequence. Instead, a long stretch of pyrimidines (approx. 18 repeats of TTC and related sequences) was found, followed by a longer stretch of purines (approx. 42 repeats of GGA or GGAAGA), and finally an A-rich sequence from which the dT-tailed primer apparently primed. Sequences similar to this have been denoted as 'Z' DNA and have been reported to decrease levels of transcription in nearby genes. The sequences are not shown in Fig. 1 due to difficulty in generating an exact sequence through the repetitive regions of high secondary structure, and because in the absence of a polyadenylation signal sequence, we cannot be confident that this region was generated from cDNA as opposed to intronic DNA. As described below, we are confident of the validity of the open reading frame and its 3' end. However, the repetitive part of the 3'UTR and the region containing the Z-DNA have to be viewed as potentially artifactual, even though they derived from legitimate RACE cDNA clones.

Gcm1 exhibits a region of extended conservation with *Drosophila gcm* (Fig. 2). The conserved block is 156 amino acids in length and overall 50% identical (63% conserved). The conserved region encompasses the cysteine-rich sequence originally described in *Drosophila gcm* (amino acids 93–148 of *Drosophila gcm*; 7 of the 8 cysteines are present in the mouse sequence) and encodes a non-conserved PEST sequence [12], but does not encode the other motifs reported including the MPVP repeat, nuclear localization sequence or homopolymeric repeats, aside from a single glutamine triplet. Thus, the protein is likely to have a short half-life and the cys-

teine-rich sequence is likely to be of biological significance for function of the extended gene family, but the other motifs are presumably not critical.

A second feature that may decrease the levels of *Gcm1* protein lies in the 5' UTR, where a small open reading frame is located. The ORF begins at nt 39 and ends with a stop codon at nt 222. In contrast, the ORF conserved with *Drosophila gcm* begins at nt 205. Since the amino-termini of *gcm* and *Gcm1* exhibit significant conservation, this finding is unlikely to result from a sequencing error or cDNA clone artifact. Such 5' UTR small ORFs have been described before, in particular in developmental genes expressed at low levels such as growth factors and homeobox genes [13,14]. In some cases, the small ORFs end one nucleotide beyond the start of the initiating codon of the biologically conserved ORF. In these instances, a start-stop-reinitiation mechanism has been proposed as a means of attenuating but not eliminating translation of the conserved ORF [15]. With respect to *Gcm1*, however, it is unlikely that the ribosomal complex would backup 17 nt; thus the conserved ORF will presumably be translated only when the first ATG is ignored, which should happen a significant fraction of the time since the ATG at nt 39 is embedded in a poor Kozak consensus sequence (TGT ATG A) as opposed to the one at the beginning of the conserved ORF, which encodes a good consensus sequence (ATC ATG G) [14].

Specific amino acid comparisons shown in Fig. 3 illustrate the high degree of similarity of *Gcm1* with the expressed human sequence tag. Conservation is observed through the carboxy-terminus, confirming the validity of the hypothetical open reading frame.

3.2. How many mammalian GCM genes exist?

The mouse chromosomal location of *Gcm1* was determined by interspecific backcross analysis using progeny derived from matings of [(C57BL/6J × *Mus spretus*)F₁ × C57BL/6J] mice. C57BL/6J and *M. spretus* DNAs were digested with several enzymes and analyzed by Southern blot hybridization for informative restriction fragment length polymorphisms (RFLPs) using a mouse 5' UTR probe. The 15.0 and 19.5 kb *Bgl* *M. spretus* RFLPs (see Section 2) were used to follow the segregation of the *Gcm1* locus in backcross mice. Surprisingly, the

```

m MELDDDFPEDKEILSWDINDVKLPQNVKTTDWFQEWPD SYVKHIYSSDDR NAQRHLSSWA MRNTNNHNSR ILKKSC LGVVVCS 83
h
d MVLNGMPITMPVMPVMPVMPSPATKS+VA+D....S..+SVGEF-.D.N+.S+G+C+L...VQSD E.++A.G.....VN.....+++ 100

m RDCSTEEGRKIYLRPAICDKARQKQQRKSCPN--CNGPLKLIPCRGHGGFVPTNFWRHDGRFIFQSKGEHDHPRPETKLEAEARRAMKKVHMASASNSLR 182
h .ARSTC.....RV.--.D.....K.....NT.PS.V..S
d AK.KLPN.AS++.....+.G.Q...RN...R.++QA....C.....+.R..NG.....T.....A.GST....L+ACGRRVRSLAVML 201

m MKGRPAAKALPAEIPSGSLPLTWSFQEGVQLPGTYSTPLIANAPQONSLNDCLSF PKSYDLGGSTEL EDPTSTLDSMKFYERCKFSSSRIYGSEEQFQP 283
h L..HRDQGLFQG.TQ.....GR [----->

m VPGTYGDYEDLQTWNKNVALGRNPDDIYYPAYPLVASWPYDYFPSQNSLEHLPQQVPSEPPAAQPGCHPLWSNPGGEFYEKVSVDLSYVP-SLTYHP 383
h <-----] .K.....T...N...V..HT.S

m PQQDPFL LTYGSPTQQQHALPGKSNRWDFDEEMACMGLDHFNEMLLNFCSLR* 436
h ..E....F..A.HPH..YS..S..SK...E...TYL....C..D...L.P..*

```

Fig. 3. Amino acid comparisons of *Drosophila gcm*, *Gcm1*, and a human *Gcm*-like expressed sequence tag. Dots indicate identical amino acids; plus signs indicate conservative changes; dashes indicate missing amino acids. Blank space indicates sequence not shown. Only the first 201 amino acids of *Drosophila gcm* are shown. Sequence information is not available for part of the human expressed sequence tag (Genbank accession no. R62635).

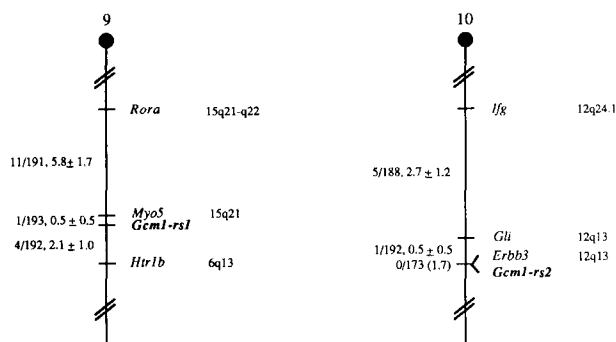


Fig. 4. Chromosomal locations of the *Gcm1* related loci in the mouse genome. *Gcm1* loci were mapped by interspecific backcross analysis. The number of recombinant N_2 animals is presented over the total number of N_2 animals typed to the left of the chromosome maps between each pair of loci. The recombination frequencies, expressed as genetic distance in centimorgans (\pm one standard error) are also shown. The upper 95% confidence limit of the recombination distance is given in parentheses when no recombinants were found between loci. Gene order was determined by minimizing the number of recombination events required to explain the allele distribution patterns. The positions of loci on human chromosomes, where known, are shown to the right of the chromosome maps. References for the map positions of most human loci can be obtained from GDB (Genome Data Base), a computerized database of human linkage information maintained by The William H. Welch Medical Library of The Johns Hopkins University (Baltimore, MD).

two fragments segregated independently; the 15.0 kb fragment mapped to the central region of mouse chromosome 9, 0.5 cM distal of *Myo5*, and defined *Gcm1-rs1* (Fig. 4). The 19.5 kb *Bgl*I fragment mapped to the distal region of mouse chromosome 10, did not recombine with *Erbb3* in 173 mice typed in common and defined *Gcm1-rs2*. The data indicate that there are at least two *Gcm1* loci in mouse and it remains to be determined which locus is the authentic homolog of the *Drosophila gcm* locus. No neurological mutants have been mapped in these regions, although there are a number of pre- and post-natal lethal factors in the corresponding intervals.

3.3. Expression analyses

Extensive expression analyses were performed using RNA in situ hybridization on embryos of various stages and neonates, and PCR was performed on cDNA mixtures prepared from varied sources. *Gcm1* could be amplified easily from placental cDNA, and with additional cycles, from E10.5 day embryos (data not shown). However, no signal was detected using RNA in situ hybridization to examine embryo sections from ages E9.5 to birth and day 4 postnatal or adult brain (data not shown). Control probes consistently detected SCG10, a marker of neural differentiation, which confirmed

that the in situ hybridization experiments were technically successful.

These results suggest two possibilities. First, it is possible that *Gcm1* is expressed in the placenta (and possibly elsewhere), but is not expressed at significant levels during embryogenesis and thus has a biological role other than in gliogenesis. Cloning and expression analysis of the second presumed *Gcm* allele may provide insight into this hypothesis. Second, a small number of developmentally critical genes have historically proven very challenging to detect by in situ hybridization, including basic FGF and activin. *Drosophila gcm* is expressed transiently, and low levels of protein are achieved through instability sequences in the RNA and PEST sequences in the protein. Since the mammalian RNA does not encode an instability sequence, it may achieve the same restricted level of expression through reduced levels of transcription.

Ectopic expression of *Gcm1* or inactivation through homologous recombination will be required to determine whether it plays a requisite and dominant role in gliogenesis.

Acknowledgements: We thank Deborah B. Householder for excellent technical assistance, and Dr. L.A. Frohman for having kindly provided mouse placental RNA. This research was supported, in part, by the National Cancer Institute, DHHS, under contract with ABL, to N.A.J., and by a National Institute of Health Grant HD29758 to M.A.F. The *Gcm1* nucleotide sequence has been deposited in the Genbank database under accession no. U59876.

References

- [1] Pfrieger, F.W. and Barres, B.A. (1995) *Cell* 83, 671–674.
- [2] Bondar, R.L. (1977) in: *Dynamic Properties of Glia Cells* (Schofeniels, E. ed.) pp. 3–11, Pergamon, New York.
- [3] Hosoya, T., Takizawa, K., Nitta, K. and Hotta, Y. (1995) *Cell* 82, 1025–1036.
- [4] Jones, B.W., Fetter, R.D., Tear, G. and Goodman, C.S. (1995) *Cell* 82, 1013–1023.
- [5] Anderson, D. (1995) *Neuron* 15, 1219–1222.
- [6] Frohman, M.A. (1994) *PCR Methods Appl.* 4, S40–S58.
- [7] Frohman, M.A., Boyle, M. and Martin, G.R. (1990) *Development* 110, 589–608.
- [8] Copeland, N.G. and Jenkins, N.A. (1991) *Trends Genet* 7, 113–118.
- [9] Jenkins, N.A., Copeland, N.G., Taylor, B.A. and Lee, B.K. (1982) *J. Virol* 43, 26–36.
- [10] Hasson, T. et al., submitted.
- [11] Green, E.L. (1981) in: *Genetics and Probability in Animal Breeding Experiments*, pp. 77–113, Oxford University Press, New York.
- [12] Rogers, S., Wells, R. and Rechsteiner, M. (1986) *Science* 234, 364–368.
- [13] Breier, G., Dressler, G.R. and Gruss, P. (1988) *EMBO J.* 7, 1329–1336.
- [14] Kozak, M. (1989) *J. Cell Biol.* 108, 229–241.
- [15] Kessel, M. and Gruss, P. (1988) *Nature* 332, 117–118.