



Holt, C. (2016) Casein and casein micelle structures, functions and diversity in 20 species. *International Dairy Journal*, 60, pp. 2-13.(doi:[10.1016/j.idairyj.2016.01.004](https://doi.org/10.1016/j.idairyj.2016.01.004))

This is the author's final accepted version.

There may be differences between this version and the published version.
You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/120662/>

Deposited on: 5 July 2016

1 **Casein and casein micelle structures, functions and**
2 **diversity in 20 species**

3 **Carl Holt**

4 **International Dairy Journal (2016)**
5 **DOI: 10.1016/j.idairyj.2016.01.004**

6 Institute of Molecular, Cell and Systems Biology, School of Life Sciences, University of
7 Glasgow, Glasgow G12 8QQ, UK.

8 Corresponding author: Carl Holt, RB413A Level B4, Joseph Black Building, University of
9 Glasgow, Glasgow G12 8QQ, Tel. +441292560158. Email carl.holt@glasgow.ac.uk

10 ORCID ID 0000-0002-2087-1546

11

12 Primary structures of caseins from 20 species, including two monotremes and two
13 marsupials, have been compared. Sequences of the mature proteins are very divergent
14 whereas variation in amino acid composition is mostly restricted to a range of disorder-
15 promoting residues. The number and size of clusters of phosphorylation sites in the caseins is
16 variable, blurring the boundaries between them. Casein polar tract sequences were found in
17 all caseins and are chiefly responsible for the weak and dynamic interactions among the
18 tangled web of peptide chains in the matrix of casein micelles. The interactions take the
19 predominant form of backbone-to-backbone contacts rather than sequence-specific side chain
20 interactions of the hydrophobic effect. It is suggested that the dynamic casein micelle matrix
21 be represented by an ensemble of interchanging structures with different types and degrees of
22 inhomogeneity, influenced by solvent quality and other environmental factors.

23

24 1. *Introduction*

25 Caseins are pleiotropic proteins with an original and continuing function in biomineralisation
26 and a better-known function as nutritional proteins. Recent research has revealed functions of
27 the casein micelle relating to the needs of the mother rather than the neonate. In particular,
28 the sequestration of calcium phosphate in a thermodynamically stable complex with caseins
29 allows milk to be stored in the mammary gland without causing benign ectopic or
30 pathological calcification in the cisternae or ducts of the gland (Holt & Carver, 2012; Holt,
31 Carver, Ecroyd, & Thorn, 2013; Thorn, Ecroyd, Carver, & Holt, 2015). The same
32 phenomenon is found in other biofluids including blood, extracellular fluid and saliva, but at
33 very much lower concentrations of the sequestering protein, and is necessary to allow the
34 easy coexistence of soft and mineralised tissues in the same organism (Holt, Sorensen, &
35 Clegg, 2009; Holt, 2013; Holt, Lenton, Nylander, Sorensen, & Teixeira, 2014). For this
36 purpose, the caseins are required to have clusters of phosphorylated residues and a disordered
37 conformation but they must also have sticky sequences so that they can form clots or a gel in
38 the stomach of the neonate. The sticky sequences in two of the individual bovine caseins
39 contain certain smaller sub-sequences that are able to form amyloid fibrils at a rate that is
40 facilitated by the disordered conformation of the whole protein (Leonil, et al., 2008; Thorn, et
41 al., 2005; Thorn, Ecroyd, Sunde, Poon, & Carver, 2008; Treweek, Thorn, Price, & Carver,
42 2011). Fibril formation is, however, suppressed in milk by an alternative association pathway
43 to produce the amorphous casein micelle. That casein genes affect three seemingly distinct
44 phenotypic traits is not anomalous because the phenomenon of pleiotropism is common and
45 forms a central concept in developmental biology (Hodgkin, 1998; Liberles, Tisdell, &
46 Grahnen, 2011; Tokuriki, Stricher, Serrano, & Tawfik, 2008; Wang, Liao, & Zhang, 2010;
47 Weinreich, Delaney, DePristo, & Hartl, 2006; Zeng & Gu, 2010). Indeed, the three functions
48 of casein genes are essential to a successful reproductive strategy to such an extent that no
49 single one should be described as the main function of caseins or the casein micelle.
50 Notwithstanding this, there is an inevitable trade-off among the functions which affects the
51 composition, structure and stability of the proteins.

52 Much of the scientific knowledge of caseins is derived from work done on farm animals,
53 particularly the cow genus *Bos*. As a result, casein research has a boocentric bias. Capuco and
54 Akers (Capuco & Akers, 2009) have pointed out that “Because no single species can provide
55 an ample and sufficient model for the physiology of another, and because the potential gain in

knowledge from comparative studies is great, the research community should not be species-centric." Accordingly, an attempt is made here to extend earlier work (Ginger & Grigor, 1999; Martin, Cebo, & Miranda, 2013) by examining sequences from 20 species. This non-random sample includes a number of species closely related to the cow and is far from being representative of extant mammals but it does include two monotremes and three marsupials. This and further work may enable a more balanced, less-boocentric, perspective on the nature of caseins to be obtained.

2. Methods

2.1 Nomenclature

Casein gene and protein nomenclature is species-specific and potentially confusing so here the four bovine genes or protein names will be used wherever possible to identify the corresponding non-bovine orthologues. These are CSN1S1 encoding α_{S1} -, CSN1S2 encoding α_{S2} -, CSN2 encoding β - and CSN3 encoding κ -casein (Lefèvre, Sharp, & Nicholas, 2010; Rijnkels, Kooiman, de Boer, & Pieper, 1997)). A recent gene duplication in rodents and some other eutherians has generated the paralogous CSN1S2A and CSN1S2B-like genes. The CSN1S2A gene is orthologous to bovine CSN1S2. Previously the sequences coded by CSN1S2B-like genes were aligned with those of CSN1S2 orthologues (Kawasaki, Lafont, & Sire, 2011) but at the protein level the differences are quite marked, justifying them being regarded as a distinct casein. In the monotremes, the CSN2A gene is orthologous to the eutherian CSN2 but the CSN2B gene appears to be a chimera in which exons 2-6 can be aligned with exons from CSN2 and the following exons with CSN1S2 (Lefevre, Sharp, & Nicholas, 2009). Thus there are 6 casein genes known to science.

2.2. Sequences and sequence alignments

The sequences used include those from the work of Kawasaki *et al.* (Kawasaki, et al., 2011). Accession codes of additional sequences, mainly for the κ -casein alignments, are specified in Supplementary Data S1 and are for the default isoforms, including the signal sequence, in the UniprotKB database on the ExPASy Bioinformatics Resource Portal (<http://www.expasy.org>). Exon boundary conservation (Martin, et al., 2013) was used as a constraint throughout. Because of extensive cryptic duplications and divergence in the long exons encoding the casein polar tract sequences of β -casein, Kawasaki *et al.* (Kawasaki, et al., 2011) found that

86 no reliable multiple alignment was obtained across eutherian, marsupial, and monotreme
87 sequences. In their work, separate alignments were made for each clade and were linked in
88 their diagram without any further attempt at inter-clade alignment. A similar problem was
89 encountered here in the alignment of 17 sequences encoded by the long exon 4 of κ -caseins.
90 The problem was alleviated in both cases, although not eliminated, by using a unit scoring
91 matrix to avoid giving undue weight to residues that are highly conserved in globular proteins
92 but less-well conserved in disordered sequences (Brown, et al., 2002; Holt & Sawyer, 1993).
93 Finally, alignments were edited manually to reduce the number of separate insertions or
94 deletions. The resulting alignments of translated sequences, including signal sequences are
95 shown in Supplementary Data S2-S6 for the 5 casein groups defined as (i) the α_{S1} -caseins
96 (CSN1S1; Data S2), (ii) orthologues of the bovine α_{S2} -casein (CSN1S2; Data S3), (iii)
97 miscellaneous other α_{S2} -type caseins (CSN1S2B-like + part monotreme CSN2B; Data S4),
98 (iv) the β -caseins (CSN2 + part of monotreme CSN2B; Data S5) and (iv) the κ -caseins
99 (CSN3; Data S6).

100 **2.3. Casein polar tracts**

101 Polar tracts are sequences found in intrinsically disordered proteins that are capable, under
102 certain circumstances, of forming a more condensed structure by interacting with themselves
103 or similar sequences. They are deficient in charged, hydrophobic, and Pro residues and
104 enriched in polar amino acids such as Asn, Gly, Gln, His, Ser, and Thr (Das, Ruff, & Pappu,
105 2015). Polar tract-type sequences that are enriched in Pro tend to form extended, water-rich
106 structures such as gels, viscous mucus and slimes rather than condensed structures (Kay,
107 Williamson, & Sudol, 2000; Williamson, 1994). Casein polar tracts are largely hydrophilic
108 sequences in which Pro and Gln are the most common residues (P,Q-rich) and they are
109 encoded by long exons. The partial exception is the long exon 4 of CSN3 (κ -casein) which
110 encodes two different types of polar tract: a P,Q,Y-rich sequence in the N-terminal half and a
111 P,S,T-rich sequence in the C-terminal half. Separating the two is an aspartate proteinase-type
112 cleavage site which generates the macropeptide in the stomach of some neonates. Here the
113 term macropeptide is applied generically to the P,S,T-rich sequences of κ -caseins whether or
114 not they are cleavage products of an aspartate proteinase

115 **2.4. Clusters of phosphorylation sites**

116 The sequences were inspected for canonical sites of phosphorylation as defined by the work
117 of Jean-Claude Mercier and his colleagues (Martin, et al., 2013; Mercier, 1981; Mercier &
118 Vilotte, 1993) and by the observed specificity of the Golgi casein kinase Fam20C
119 (Tagliabronci, et al., 2012). The canonical kinase recognition sequence is -S-X-E/pS-, where
120 X is any residue and pS is phosphoserine. Actual degrees of phosphorylation may be less than
121 predicted due to incomplete phosphorylation of the available sites but may be larger because
122 of phosphorylation at non-canonical sites such as -T-X-E/pS/D- or -S-X-D-. A cluster of n
123 phosphorylation sites results from complete phosphorylation of sequences such as -S-(S)_{n-1}-
124 E-E or -S-X-(S)_{n-1}-E-E where there is a single initial kinase recognition site. In caseins, the
125 initial kinase recognition site of a cluster of sites is often formed by exon splicing (Mercier,
126 1981; Mercier & Vilotte, 1993). An exception is the potential phosphate cluster in platypus κ -
127 casein which is found near the N-terminus of the casein polar tract sequence encoded by exon
128 4 (-S-S-E-E-S-E-E-). It is therefore not formed by exon splicing and contains two initial
129 kinase recognition sites. A smaller cluster is found at the same position in the echidna κ -
130 casein.

131 A peptide containing a phosphate cluster may be capable of sequestering amorphous calcium
132 phosphate to form a thermodynamically stable complex. The requisite number and nature of
133 phosphorylated residues in such a cluster is not well defined experimentally but all known
134 examples have had three or more pS residues close together with a single initial kinase
135 recognition site (Clegg & Holt, 2009; Holt, Wahlgren, & Drakenberg, 1996; Holt, Timmins,
136 Errington, & Leaver, 1998; Holt, et al., 2009; Little & Holt, 2004).

137 **2.5. Other sequence analysis tools**

138 Amyloid-forming zipper sequences were predicted by the method of Goldschmidt *et al.*
139 (Goldschmidt, Teng, Riek, & Eisenberg, 2010) and edited to remove predictions that
140 conflicted with likely post translational modifications of phosphorylation and disulphide bond
141 formation (Holt & Carver, 2012).

142 Amino acid composition and the mole fractions of positively or negatively charged residues
143 were calculated using the ProtParam web tool (Gasteiger, et al., 2005)
144 (<http://web.expasy.org/protparam/>). Hydropathy values were calculated from the

145 hydrophobicity scale of Kyte and Doolittle (Kyte & Doolittle, 1982) and renormalized to a
146 scale from zero (Arg) to one (Ile) (Uversky, Gillespie, & Fink, 2000). The disorder
147 propensity scale was taken from the TOP-IDP analysis (Campen, et al., 2008).

148 **3. Results**

149 **3.1. Structure of casein genes**

150 The phylogenetic studies of Kawasaki *et al.* (Kawasaki & Weiss, 2003; Kawasaki, et al.,
151 2011) showed that all caseins are members of a group of secreted, calcium (phosphate)-
152 binding proteins. They are formed from short and long exons separated by phase zero introns
153 and the first and last exons of casein genes are totally untranslated. Exon 2 encodes a signal
154 sequence which is necessary for secretion. All casein genes evolved from the bone and tooth
155 gene ODAM which encodes the ODontogenic AMeloblast-associated protein but the
156 calcium-sensitive casein genes took a slightly different subsequent path to the κ -casein gene.
157 Thorn *et al.* (Thorn, et al., 2015) attempted to capture the essential differences between the
158 two classes of casein genes by means of a Bauplan for each. The relationship of the Bauplan
159 to some actual CSN2 gene structures is illustrated in Figure 1

160 **3.2. Casein composition**

161 The composition of whole casein is known to a reasonable degree of completeness in only a
162 handful of eutherian species (Holt, et al., 2013; Martin, et al., 2013). Even within this small
163 group, the proportions of the individual caseins are widely variable. For all the caseins except
164 κ -casein, the proportion may be zero in certain individuals of a species, or in all members of
165 a species or group of related species or in whole lineages. This is in spite of the evolution of
166 the four orthologues CSN1S1, CSN1S2, CSN2 and CSN3 before the radiation of mammals
167 (Lefèvre, et al., 2010; Oftedal, 2012). From the limited data to hand at the present time, it
168 appears that casein micelles are formed from a mixture of between three and five types of
169 casein, depending on the species. The nutritional function of the casein micelle does not
170 provide a compelling explanation for why a mixture of casein types is preferred since the
171 caseins have rather similar amino acid compositions and none is rich in essential amino acids
172 (Hamraeus & Lonnerdal, 2003). By contrast, an explanation in terms of the suppression of
173 amyloid formation by individual caseins (Holt & Carver, 2012; Holt, et al., 2013) has
174 experimental support (Thorn, et al., 2005; Thorn, et al., 2008; Treweek, et al., 2011). The

175 experiments show that in the bovine caseins the formation of amyloid fibrils by any
176 individual casein is suppressed by all the other casein types in the mixture, each of which acts
177 as a molecular chaperone. Thus, a mixture of different types of casein is preferable to one,
178 amyloid-forming, casein. What is still not clear is what determines the number of components
179 in the mixture and their proportions. Also, why is it not possible to naturally select caseins
180 without amyloid-forming tendencies? The motivation for the present work is that studies of
181 non-bovine caseins from as wide a range of species as possible should help to solve these and
182 related problems.

183 The mean mole % and standard deviation of residues in the mature translated but
184 unphosphorylated sequences of the CSN1S1, CSN1S2, CSN2 and CSN3 orthologues are
185 shown in Figure 2(a)-(d). Excluded from this analysis are the eutherian CSN1S2B-like
186 sequences and the monotreme CSN2B sequences because there are few examples of these
187 from which to calculate reliable means with small standard deviations.

188 The most common residues are very largely drawn from the group of non-essential amino
189 acids with a high disorder propensity (Figure 2e) and low hydropathy (Figure 2f).
190 Hydrophobic residues, normally associated with the core of globular proteins, occur
191 relatively infrequently in all members of the sampled group. Only one essential, hydrophilic,
192 amino acid residue occurs at high frequency, and only in the α_{S2} -caseins and this is Lys,
193 which has a high disorder propensity. Only one hydrophobic essential amino acid with a low
194 disorder propensity, Leu, occurs at high (in the β -caseins) or moderate (in the α_{S1} -caseins)
195 frequency. The essential amino acids Thr and Val are found at high-to-moderate frequency in
196 κ -caseins and Val also occurs at moderate frequency in β -caseins.

197 There are relatively few residues that can be used to discriminate among the casein
198 orthologues. These can be identified in the multiple bar chart of mean compositions shown in
199 Figure 3a. Lys in α_{S2} -, Leu in β - and Thr and Ala in κ -casein occur more frequently than in
200 the other orthologues. Thr in α_{S1} - and Tyr in β -caseins occur less commonly, on average,
201 than in the other orthologues. Figure 3b shows a 3-D scatter plot of the individual species
202 using the orthogonal axes of the mole % Lys (x), mole % Leu (y) and mole % Thr (z).
203 Inspection of this projection and of others showed that the orthologous groups were
204 reasonably well separated using only these three dimensions.

205 At this point, any explanations for the compositional differences among the caseins are still
206 highly speculative. The higher frequency of Thr residues in the κ -caseins may be due in part
207 to a need for actual or potential sites for *O*-glycosylation in the macropeptide sequences. The
208 Leu and Val residues in β -caseins occur in the casein polar tract sequences and may serve to
209 supplement the backbone-to-backbone interactions that predominate in these regions with
210 hydrophobic interactions of the side chains, either stabilising intermolecular binding or
211 reinforcing a range of preferred intramolecular conformations. The higher frequency of Lys
212 residues in the α_{S2} -casein A orthologues may be to reduce the net negative charge arising
213 from phosphorylated residues but the positively charged residues are concentrated towards
214 the C-terminus and hence that region could help to stabilise a range of dynamic
215 conformations in which it interacts electrostatically with more central regions of net negative
216 charge.

217 **3.3. Multiple sequence alignment**

218 The alignments shown in Supplementary Data S2-S6 are aligned in boxes corresponding to
219 translated exon sequences such that all entries in the same box are considered to be
220 orthologous (Kawasaki, et al., 2011). Fully conserved residues (identities) are those that
221 appear in every species at a given position in the alignment. There were very few (63)
222 positions of complete conservation in the 5 alignments and 24 of these were in the
223 CSN1S2B-like group where diversity is low because there are only a few aligned sequences.
224 Of the remaining 39 identities, most (22) were in sequences encoded by exon 2, 7 were in
225 polar tracts and 10 were in sequences encoded by short exons, such as phosphate clusters (7).
226 The signal sequences of the κ -caseins appear to be less well conserved than in other
227 caseins. Residues such as Cys and Trp are normally better conserved in globular proteins than
228 Ser or Glu but in the caseins the reverse is true, as was previously noted in an alignment of a
229 smaller number of eutherian caseins (Holt & Sawyer, 1993). Cys residues in caseins differ
230 from most Cys residues in globular proteins in that they are almost exclusively involved in
231 intermolecular disulphide bridges but their role, if any, in the maintenance of casein micelle
232 structure remains an enigma(Thorn, et al., 2015).

233 Fully conserved exons are those that are found in every species in an alignment. The numbers
234 of fully conserved translated exons were 5 (13%), 9 (47%), 7 (41%), 2 (20%) and 3 (100%)
235 in the alignments of Supplementary Data S2-S6, respectively. Here, the numbers in

236 parentheses are the fraction of fully conserved exons as a percentage of the total number of
237 exons in the alignment. Not surprisingly, the number of conserved positions and exons is
238 considerably higher when only eutherian species are considered. For example, among the
239 eutherian β -caseins, 5 of the 6 translated exons are fully conserved compared to only two
240 among the 9 translated exons of all the β -caseins and the β -casein-like part of monotreme
241 CSN2B. The number of conserved residues and exons depend, to a degree, on the
242 assumptions and methods used in the alignment and the particular sequences chosen. A more
243 diverse and fully representative group of mammalian casein genes would almost certainly
244 have produced even fewer examples of complete conservation but the chosen sample does
245 demonstrate clearly that apart from the signal peptides of the calcium-sensitive caseins,
246 casein sequences are very variable. This conclusion reinforces previously expressed views
247 based on smaller groups of eutherian species (Martin, et al., 2013; Mercier & Villette, 1993)
248 that the mature caseins show very low levels of sequence conservation.

249 **3.4. Occurrence and variation of phosphorylation site clusters**

250 The clustering of phosphorylation sites has been analysed by dividing the sequences into 4
251 groups. These groups are the sequences from Supplementary Data S2 (α_{S1} -caseins),
252 Supplementary Data S3 combined with Supplementary Data S4 (α_{S2} -type caseins),
253 Supplementary Data S5 (β -caseins) and Supplementary Data S6 (κ -caseins). A histogram
254 showing the frequency of occurrence of different cluster sizes in the individual casein groups
255 and for all caseins is shown in Figures 4a and 4b, respectively.

256 Whereas single sites were the most common cluster size in κ -caseins (Fig. 4a) and the most
257 common in the aggregate of all calcium-sensitive caseins (Figure 4b), the most common size
258 in the individual calcium-sensitive caseins was 3 (α_{S2} -like), 4 (β -) or 5 (α_{S1} -) and the largest
259 size was 7 in guinea pig, platypus and echidna α_{S1} -caseins. Figures 4c and 4d show the
260 number of clusters of size $n \geq 2$ or $n \geq 3$ in individual sequences. The number varied from
261 zero in most κ -caseins to 3 ($n \geq 3$) in some of the α_S -caseins. For a cluster size of 2 or greater
262 the maximum number of clusters in a sequence was 5 in guinea pig α_{S2} -casein.

263 A minimum cluster size of $n = 3$ is required for their cross-linking action in casein micelles
264 according to the findings of Aoki and co-workers (Aoki, Umeda, & Kako, 1992; Umeda &

265 Aoki, 2002). In the formation of calcium phosphate nanoclusters using pure phosphopeptides,
266 thermodynamically stable complexes have only been observed if $n \geq 3$ (Holt, et al., 1996;
267 Holt, et al., 1998; Little & Holt, 2004). Nevertheless, in mixtures of caseins or casein and
268 osteopontin phosphopeptides, smaller clusters of $n = 2$ appear to enter the sequestering shell
269 of peptides provided larger cluster sizes are also present (Clegg & Holt, 2009; Holt, et al.,
270 2009). The theory of amorphous calcium phosphate sequestration by phosphopeptides has
271 been developed using the term “phosphate centre”, defined as a short sequence containing
272 three or more phosphorylated residues. Ideally, this simplifying concept, although it has been
273 very useful, will be replaced one day by a theory that accounts explicitly for the effects of the
274 phosphate cluster size distribution. In eight eutherian species for which casein composition is
275 well-enough established, the average number of phosphate centres per average mole of
276 caseins has been calculated and shown to vary only within close bounds (Holt & Carver,
277 2012). If confirmed for a broader range of species it would indicate that the main factor in the
278 interspecific variation of sequestered calcium and phosphate concentrations is the
279 concentration of casein rather than its composition or the number and distribution of potential
280 sites of phosphorylation.

281 When the data in Figure 4 are examined as a whole, it is clear that while there are differences
282 in the average cluster size and number of clusters among the groups, the distributions overlap
283 considerably, blurring the boundaries between caseins in their ability to sequester amorphous
284 calcium phosphate.

285 **3.5. Conservation of casein polar tracts**

286 All the sequences in the sample contain at least one casein polar tract (coloured green in
287 Supplementary Data S2-S6) and in the α_{S2} -casein orthologues and in the rabbit and rat α_{S2} -
288 casein B sequences there are two because of an intragenic duplication (Stewart, et al., 1987).
289 This level of conservation is second only to the conservation of the signal peptide sequences
290 in caseins. The variation in the lengths of the polar tracts has been investigated by dividing
291 the sequences into the same 4 groups used in the phosphorylation site cluster analysis.
292 Among the four groups there are clear differences in the average lengths of the casein polar
293 tract sequences with only limited overlaps in the number distribution histograms (Figure 5).
294 For example, the average number of residues in polar tract sequences in the 4 groups are
295 57.75 ± 12.76 in α_{S1} -caseins), 40.93 ± 2.30 and 40.79 ± 2.08 in the α_{S2} -like caseins, 88.68 ± 5.20

296 in the κ -caseins and 170.43 ± 28.62 in the β -caseins. Thus, the β -casein orthologues contain up
297 to approximately twice as many residues in polar tract sequences as in all other casein types.

298 **3.6. Occurrence of amyloid zipper sequences in caseins**

299 There is a limited amount of evidence that amyloid protofibrils are part of the structure of
300 bovine casein micelles (Lencki, 2007). Casein sequences were therefore examined to see
301 whether predicted amyloid zipper sequences are conserved, as they would be if the
302 protofibrils are important for casein micelle formation. Amyloid zipper predictions are
303 abundant in the casein polar tract sequences (Holt & Carver, 2012) but because of the
304 uncertain alignment of the casein polar tracts in the β -caseins, the sequences of choice to test
305 the hypothesis are in the shorter polar tracts of κ - and α_S - caseins. Amyloid zipper predictions
306 have been validated in experimental studies with other proteins (Goldschmidt, et al., 2010;
307 Nelson, et al., 2005; Rodriguez, et al., 2015; Sawaya, et al., 2007; Thompson, et al., 2006;
308 Wiltzius, et al., 2008). They were found in the regions of bovine α_{S2} - and κ -caseins where
309 amyloid formation has been demonstrated experimentally (Ecroyd, et al., 2008; Ecroyd,
310 Thorn, Liu, & Carver, 2010; H. M. Farrell, Cooke, Wickham, Piotrowski, & Hoagland, 2003;
311 Niewold, Murphy, Hulskamp-Koch, Tooten, & Gruys, 1999). Nevertheless, neither these nor
312 any other predicted amyloid zipper sequences were fully conserved in the aligned casein
313 sequences. This is illustrated in Figure 6 for the κ -caseins where the predicted amyloid zipper
314 sequences are highlighted blue in the multiple sequence alignment.

315 The complete absence of experimental evidence of amyloid fibril formation in non-bovine
316 caseins is not evidence of absence. It is important to gather more data to help explain why
317 potentially cytotoxic, amyloid-forming sequences, are tolerated in a vital food.

318 **3.7. Condensation of casein polar tracts and casein micelle structure**

319 Casein micelles have been found in all milks so far examined but a striking conclusion from
320 interspecific studies (Martin, et al., 2013) is that they can be made in a large number of
321 distinct ways using a mixture of different caseins in variable proportions.

322 The quality of the solvent and its structure in the solvation sheath around the backbone are
323 factors that are thought to be important in the condensation of polar tracts (Das, et al., 2015;
324 R. V. Pappu, Srinivasan, & Rose, 2000; Rohit V. Pappu, Wang, Vitalis, & Crick, 2008; Tran,

325 Mao, & Pappu, 2008) and for the conformational preferences of intrinsically disordered
326 proteins, especially for transitions within the conformational space enclosing the poly-*L*-
327 proline helix and the more extended β -strand structure (Ilawe, Raeber, Schweitzer-Stenner,
328 Toal, & Wong, 2015; Meral, Toal, Schweitzer-Stenner, & Urbanc, 2015). Thus, changes in
329 water structure are important in both polar tract interactions and the hydrophobic effect
330 between side chains.

331 The arguments for using the term “casein polar tract” rather than “hydrophobic tail” to
332 describe the P,Q-rich casein sequences were set out previously (Holt, et al., 2013). The
333 finding that caseins act on each other to control amyloid fibril formation and micelle size and
334 act on a broad range of denatured globular proteins to limit their aggregation is one of the
335 strongest arguments that the interactions have a low sequence specificity and are therefore
336 mainly backbone-to-backbone.

337 Casein polar tracts contain many Pro residues which favour the PP-II conformation, prevent
338 β -strands from forming and inhibit amyloid structures. In salivary Pro-rich proteins, the side
339 chain of Pro can stack with polyphenolic molecules such as tannins to form an insoluble
340 complex (Bennick, 2002; Charlton, Haslam, & Williamson, 2002; Luck, et al., 1994; Murray,
341 Williamson, Lilley, & Haslam, 1994; Pascal, Pate, Cheynier, & Delsuc, 2009; Williamson,
342 1994). Stacking (also called π - π stacking) refers to attractive, noncovalent interactions
343 between aromatic rings, since they contain π bonding electrons. Despite intense experimental
344 and theoretical interest, there is no unified description of the factors that contribute to
345 stacking interactions. The interaction of the major polycyclic flavan-3-ol from green tea with
346 the salivary proline-rich protein IB5 appears to favour the PP-II conformation (Pascal, et al.,
347 2009). In studies by NMR employing time-averaged nuclear Overhauser measurements of a
348 model peptide Q-G-R-P-P-Q-G with the polyphenol (-)-epigallocatechin gallate (Charlton, et
349 al., 2002), a range of possible conformations was generated in which there were stacking
350 interactions of the Pro residues with the A and C rings of the polyphenol. The structures
351 appear to show how a precipitate can grow from the complex by further stacking interactions.
352 In the interaction of casein micelles with the same polyphenol, it was found that up to a
353 million molecules of the tannin could be incorporated in the micelles but without forming a
354 precipitate (Shukla, Narayanan, & Zanchi, 2009).

355 Although Pro residues exhibit the stacking interaction with polycyclic phenols it is unclear
356 whether the interaction has a significant contribution from the hydrophobic effect. Pro
357 residues are not classified as hydrophobic because they hold tightly onto their solvating water
358 molecules, partly because the backbone carbonyl is a good hydrogen bond acceptor (Theillet,
359 et al., 2013). As a result, the presence of Pro residues in polar tracts restricts the condensation
360 process so that they form water-rich structures such as gels, mucus and slimes (Williamson,
361 1994). The water-rich matrix of casein micelles may not be fully homogeneous. Between
362 most of the Pro residues in caseins are short, conventional polar tract sequences, on average
363 comprising 5-6 residues , that could interact with other, similar, sequences and form more
364 compact or condensed substructures within the matrix. To explain the amplitude of a feature
365 in the small-angle scattering of casein micelles, de Kruif et al. (de Kruif, Huppertz, Urban, &
366 Petukhov, 2012) proposed that condensed protein structures on a scale of 2 nm size are
367 present in the matrix as a result of hydrophobic interactions. Hydrophobic interactions are not
368 needed for the explanation, however. In the interaction of casein polar tracts through
369 backbone-to-backbone interactions there is nevertheless a dependence on solvent quality
370 through the amino acid composition (not sequence) of the tract. Under any given condition of
371 solvent quality, fluctuations about the mean density of the matrix in casein micelles could
372 occur through alternative and nearly equivalent polar tract interactions. Other studies have
373 provided evidence of voids or channels in the protein matrix and of distortions of micelle
374 shape at a surface or in ice (Bouchoux, Gésan-Guiziou, Pérez, & Cabane, 2010; Bouchoux, et
375 al., 2015; Dagleish, Spagnuolo, & Goff, 2004; Gebhardt & Kulozik, 2014; Ouanezar,
376 Guyomarc'h, & Bouchoux, 2012; Trejo, Dokland, Jurat-Fuentes, & Harte, 2011). Many older
377 electron microscopy techniques produced images of larger-scale substructures called
378 submicelles and a debate continues on the extent to which these submicelles are artefacts of
379 the sample preparation methods (Farrell, Malin, Brown, & Qi, 2006; McMahon & McManus,
380 1998; McMahon & Oommen, 2012). Drying changes the small-angle X-ray scattering of
381 casein micelles, including effects on the length scale attributed to submicelles or the average
382 distance between nanoclusters (Mata, Udabage, & Gilbert, 2011).

383 Irrespective of whether any particular structure is an artefact, as some undoubtedly are, a
384 generalisation can be made from these various observations which is that the structure of the
385 casein micelle is both fragile and dynamic and its internal structure can be perturbed in a

386 variety of different ways so that internal fluctuations in matrix density become more or less
387 important.

388 An intrinsically disordered protein can be described by its average size and shape. However,
389 such proteins are dynamic and explore a huge number of alternative conformations, some of
390 which are preferred but most are transient. The description of all possible conformations is
391 impracticable but an improvement over the description of the average size and shape is the
392 use of an ensemble of structures derived from scattering and spectroscopic measurements
393 (Bernado, Bertoncini, Griesinger, Zweckstetter, & Blackledge, 2005; Bernado, Blanchard, et
394 al., 2005; Jensen, Salmon, Nodet, & Blackledge, 2010). The ensemble describes both the
395 average structure and of excursions from the average structure experienced by a single
396 molecule over time or, equivalently, by a population of molecules at an instant of time. The
397 ensemble is a collection of static structures that conveys the dynamic or mutable nature of the
398 molecule they represent. The weighting given to a particular member of the ensemble can be
399 considered to be the size of the subpopulation having that conformation. Thus, the effect of a
400 change of environment such as solvent quality can be represented by shifts in the population
401 among ensemble members.

402 The dynamic nature of casein micelles means that at any one moment a micelle can have a
403 substructure which departs in some way from the average. The ensemble hypothesis for
404 describing native casein micelles is that the total range of substructures can be represented by
405 a much smaller ensemble of substructures such that the average substructure and the
406 distribution of alternative substructures is obtained as a weighted sum over all members of
407 the ensemble. In addition to a distribution of substructures, casein micelles exist as a dynamic
408 distribution of sizes, also responsive to environmental change, although the rate of change
409 may be slow (Dewan, Chudgar, Mead, Bloomfield.V, & Morr, 1974; Huppertz, Kelly, & de
410 Kruif, 2006; Jackson & McGillivray, 2011; Lin, Dewan, Bloomfield, & Morr, 1971)

411 Four members are suggested as an ensemble of alternative casein micelle substructures
412 (Figure 7a). These are all nanocluster models having various types or degrees of
413 disproportionation of the protein matrix density. These are (7b) a more-or-less homogeneous
414 matrix, (7c) a matrix with void spaces, (7d) a matrix with condensed protein structures and

415 (7e) a matrix with both void spaces and condensed protein structures. This basis set can be
416 modified or added to in the future as the need arises.

417 Experiments on native casein micelles diluted with their own ultrafiltrate and studied by
418 small-angle X-ray or neutron scattering suggest that the protein matrix is relatively, though
419 not completely, homogeneous (de Kruif, et al., 2012; de Kruif, 2014; Holt, de Kruif, Tuinier,
420 & Timmins, 2003; Ingham, et al., 2015; Marchin, Putaux, Pignon, & Léonil, 2007; Shukla, et
421 al., 2009). Nevertheless, it very readily disproportionates to form less homogeneous
422 structures under a wide range of conditions including heat treatment, drying, freezing, the
423 reduction of water activity, adsorption at a surface, filtration forces, addition of ethanol or
424 other poor solvent and addition of polyphenols. The fact that disproportionated matrices are
425 readily observed is an indication that such structures are already present in an ensemble of
426 native states but perhaps are not highly populated.

427 ***3.8. Pleiotropy of caseins***

428 The first casein evolved in some stem amniote, before the great divergence into the sauropsid
429 and synapsid lineages, about 50 million years before the emergence of lactation in the
430 mammal-like reptiles (Kawasaki, et al., 2011). Caseins are therefore pleiotropic proteins in
431 which an antecedent function, probably in the control of some aspect of biomineralisation, is
432 closely related to one of the current functions of sequestering amorphous calcium phosphate
433 in the form of thermodynamically stable nanoclusters, but unrelated to an additional current
434 function, namely, neonatal nutrition. The additional function has required a great increase in
435 expression level and introduced conflicting pressures on the composition and structure of the
436 caseins. In biofluids other than milk, such as blood, extracellular fluid, saliva and urine, the
437 sequestering function of phosphoproteins is exercised at concentrations typically three orders
438 of magnitude smaller than in milk (Holt, et al., 2014). Such an increase in concentration has
439 brought with it problems that were of minor or of no importance in the antecedent function of
440 casein.

441 An unfolded, intrinsically disordered, conformation appears to be a characteristic functional
442 feature of many of the phosphoproteins involved in the control of biomineralisation (Holt, et
443 al., 2009; Kalmar, Homola, Varga, & Tompa, 2012) and was recognised to be advantageous
444 for the sequestration of amorphous calcium phosphate by caseins and other phosphoproteins

such as osteopontin (Holt & Sawyer, 1993; Holt, 2013). However the subset of amino acids that favour disordered conformations do not overlap well with the subset of essential amino acids, as demonstrated in Figure 2e. Indeed, the amino acid composition of caseins, with few exceptions, is commensurate with the sequestering function and the delivery of high concentrations of calcium and phosphate for bone growth but it does not deliver high concentrations of the essential amino acids to the neonate.

The effect of pleiotropy on the evolution of proteins has received a great deal of theoretical and experimental attention (Delgado, et al., 2001; Hodgkin, 1998; Liberles, et al., 2011; Wang, et al., 2010). Because most mutational changes to a globular protein cause it to lose stability, the translational selection hypothesis is that protein evolution rate is controlled by protein stability and the need to avoid the formation of cytotoxic misfolding products such as amyloid (Zeng & Gu, 2010). For example, in the adaptation of a bacterial β -lactamase to a new antibiotic, a total of 5 residue substitutions produced a substantial increase in resistance to the drug but the experimental generation of the $5! = 120$ trajectories produced mostly destabilised, amyloid-forming intermediates (Weinreich, et al., 2006). In a recent impressive study of the bacterial signalling kinase, PhoQ, a total of 160,000 mutations were generated at 5 residue positions affecting the binding to its substrate, PhoP (Podgornaia & Laub, 2015). Of these 1659 were functional in that they would still recognise PhoP. Podgornaia and Laub (2015) suggest that not all the functional variants are found in nature because of further context-sensitive constraints. Thus, trajectories that somehow sustain neutral or maladaptive intermediates on the way to an adaptive new function must exist, even for globular proteins.

The problem of maladaptive intermediates is no less severe for pleiotropic intrinsically disordered proteins than for globular proteins because the generation of amyloid-forming or other types of cytotoxic sequences cannot be hidden in a hydrophobic core. The problem of amyloid fibril formation by caseins is made more severe by the pressure to increase the protein concentration to fulfil the nutritional function. In considering the evolution of caseins and milk, Holt and Carver (2012) proposed that pleiotropy favoured an increase in casein gene complexity through an epistatic mechanism. The increased complexity at the casein locus took the form of an increased number of similar casein genes. Thus, when a number of casein gene products are expressed together, each has been shown to act as a molecular chaperone to inhibit the formation of amyloid fibrils by all of the other caseins. The

476 alternative pathway of aggregation produces an amorphous casein micelle rather than highly
477 ordered amyloid fibrils. In consequence, the total casein concentration could be increased
478 safely from μ M to mM concentrations (Holt, et al., 2013) without causing amyloid to form.
479 In other words, the casein micelle has sustained a trajectory of evolutionary change producing
480 increased levels of expression and mutational changes in caseins that would otherwise be
481 maladaptive.

482 Among existing orthologues studied here there is none that has an amino acid composition
483 rich in all the essential amino acids, but as the studies of Podgornaia and Laub (2015)
484 suggest, the orthologues of caseins that exist today and all that have ever existed may still
485 only be a fraction of the total number of functional sequences. Among the latter may be
486 sequences that are better suited, in terms of essential amino acid composition, to the current
487 nutritional role of caseins.

488 Negative pleiotropy has been invoked in relation to the adaptation of signalling proteins and
489 domains such as CH2, SH3 and WW (Liberles, et al., 2011; Zarrinpar, Park, & Lim, 2003).
490 Such domains can respond in a signalling cascade to a range of ligands (Uversky, Oldfield, &
491 Dunker, 2005) but in a positive pleiotropic adaption the danger is that they will also acquire
492 the ability to bind a new ligand producing a potentially cytotoxic response. Negative
493 pleiotropy is a process in which undesirable adaptations are eliminated while new functions
494 are acquired. The ability of signalling domains to respond to different ligands is sometimes
495 called promiscuous ligand binding. Promiscuous binding of target proteins is an essential
496 attribute of molecular chaperones in preventing misfolded proteins from aggregating or
497 forming amyloid fibrils (Barral, Broadley, Schaffar, & Hartl, 2004; Ecroyd & Carver, 2009;
498 Westerheide, Raynes, Powell, Xue, & Uversky, 2012). Promiscuous interactions have been
499 considered important in the pleiotropy of viral proteins (Habchi & Longhi, 2012). The
500 interaction of co-secreted caseins leading to the formation of the casein micelle also reduces
501 the possibility of very high concentrations of caseins binding promiscuously to unintended
502 targets in the intracellular secretion pathway, the cisterns and ducts of the mammary gland
503 and the stomach of the neonate.

504 In summary, pleiotropy provides us with a potential explanation for the manifest diversity in
505 casein sequences and an increase over evolutionary time in the complexity of the casein gene

506 locus. Co-secretion of a mixture of caseins to produce the casein micelle reduced their
507 individual ability to interact promiscuously with non-casein proteins and has neutralised what
508 would otherwise be maladaptive mutations on a trajectory towards the newly acquired
509 nutritional function.

510 **4. Conclusions**

511 Caseins are pleiotropic proteins in which the antecedent function in the control of some
512 aspect of biomineralisation is related to their current function in neonatal nutrition where they
513 sequester amorphous calcium phosphate. Potentially pathological consequences of a very
514 large increase in expression level have been neutralised by increasing the number of co-
515 expressed casein genes so that they bind to each other and form the casein micelle. The
516 content of essential amino acids has remained at a low level, probably because higher levels
517 would be incompatible with the unfolded conformation needed for caseins to form a
518 thermodynamically stable complex with amorphous calcium phosphate. The casein micelle is
519 a fragile and dynamic structure which can therefore be represented better by an ensemble of
520 interconverting states than by an average state.

521 **5. Acknowledgements**

522 John Carver and David Thorn of the Australian National University , Heath Ecroyd of the
523 University of Wollongong and Jared Raynes of CSIRO, Werribee, Australia all played
524 important roles in the development of the ideas in this paper.

525 **6. References**

526 **References**

- 527 Aoki, T., Umeda, T., & Kako, Y. (1992). The least number of phosphate groups for cross-
528 linking of casein by colloidal calcium phosphate. *Journal of Dairy Science*, 75, 971-
529 975.
- 530 Barral, J. M., Broadley, S. A., Schaffar, G., & Hartl, F. U. (2004). Roles of molecular
531 chaperones in protein misfolding diseases. *Seminars in Cell & Developmental
532 Biology*, 15, 17-29.
- 533 Bennick, A. (2002). Interaction of plant polyphenols with salivary proteins. *Critical Reviews
534 in Oral Biology & Medicine*, 13, 184-196.
- 535 Bernado, P., Bertoncini, C. W., Griesinger, C., Zweckstetter, M., & Blackledge, M. (2005).
536 Defining long-range order and local disorder in native alpha-synuclein using residual
537 dipolar couplings. *Journal of the American Chemical Society*, 127, 17968-17969.

- 538 Bernado, P., Blanchard, L., Timmins, P., Marion, D., Ruigrok, R. W. H., & Blackledge, M.
539 (2005). A structural model for unfolded proteins from residual dipolar couplings and
540 small-angle x-ray scattering. *Proceedings of the National Academy of Sciences of the*
541 *United States of America*, 102, 17002-17007.
- 542 Bouchoux, A., Gésan-Guiziou, G., Pérez, J., & Cabane, B. (2010). How to Squeeze a Sponge:
543 Casein Micelles under Osmotic Stress, a SAXS Study. *Biophysical Journal*, 99, 3754-
544 3762.
- 545 Bouchoux, A., Ventureira, J., Gesan-Guiziou, G., Garnier-Lambrouin, F., Qu, P., Pasquier,
546 C., Pezennec, S., Schweins, R., & Cabane, B. (2015). Structural heterogeneity of milk
547 casein micelles: A SANS contrast variation study. *Soft Matter*, 11, 389-399.
- 548 Brown, C. J., Takayama, S., Campen, A. M., Vise, P., Marshall, T. W., Oldfield, C. J.,
549 Williams, C. J., & Dunker, A. K. (2002). Evolutionary rate heterogeneity in proteins
550 with long disordered regions. *Journal of Molecular Evolution*, 55, 104-110.
- 551 Campen, A., Williams, R. M., Brown, C. J., Meng, J. W., Uversky, V. N., & Dunker, A. K.
552 (2008). TOP-IDP-scale: A new amino acid scale measuring propensity for intrinsic
553 disorder. *Protein and Peptide Letters*, 15, 956-963.
- 554 Capucco, A. V., & Akers, R. M. (2009). The origin and evolution of lactation. *Journal of*
555 *Biology (London)*, 8, Article No.: 37.
- 556 Charlton, A. J., Haslam, E., & Williamson, M. P. (2002). Multiple conformations of the
557 proline-rich protein/epigallocatechin gallate complex determined by time-averaged
558 nuclear Overhauser effects. *Journal of the American Chemical Society*, 124, 9899-
559 9905.
- 560 Clegg, R. A., & Holt, C. (2009). An *E. coli* over-expression system for multiply-
561 phosphorylated proteins and its use in a study of calcium phosphate sequestration by
562 novel recombinant phosphopeptides. *Protein Expression and Purification*, 67, 23-34.
- 563 Dalgleish, D. G., Spagnuolo, P. A., & Goff, H. D. (2004). A possible structure of the casein
564 micelle based on high-resolution field-emission scanning electron microscopy.
565 *International Dairy Journal*, 14, 1025-1031.
- 566 Das, R. K., Ruff, K. M., & Pappu, R. V. (2015). Relating sequence encoded information to
567 form and function of intrinsically disordered proteins. *Current Opinion in Structural*
568 *Biology*, 32, 102-112.
- 569 de Kruif, C. G., Huppertz, T., Urban, V. S., & Petukhov, A. V. (2012). Casein micelles and
570 their internal structure. *Advances in Colloid and Interface Science*, 171, 36-52.
- 571 de Kruif, C. G. (2014). The structure of casein micelles: a review of small-angle scattering
572 data. *Journal of Applied Crystallography*, 47, 1479-1489.
- 573 Delgado, S., Casane, D., Bonnaud, L., Laurin, M., Sire, J.-Y., & Girondot, M. (2001).
574 Molecular Evidence for Precambrian Origin of Amelogenin, the Major Protein of
575 Vertebrate Enamel. *Molecular Biology and Evolution*, 18, 2146-2153.
- 576 Dewan, R. K., Chudgar, A., Mead, R., Bloomfield, V. A., & Morr, C. V. (1974). Molecular
577 weight and size distribution of bovine milk casein micelles. *Biochimica Et Biophysica*
578 *Acta*, 342, 313-321.

- 579 Ecroyd, H., Koudelka, T., Thorn, D. C., Williams, D. M., Devlin, G., Hoffmann, P., &
580 Carver, J. A. (2008). Dissociation from the oligomeric state is the rate-limiting step in
581 fibril formation by kappa-casein. *Journal of Biological Chemistry*, 283, 9012-9022.
- 582 Ecroyd, H., & Carver, J. A. (2009). Crystallin proteins and amyloid fibrils. *Cellular and*
583 *Molecular Life Sciences*, 66, 62-81.
- 584 Ecroyd, H., Thorn, D. C., Liu, Y., & Carver, J. A. (2010). The dissociated form of κ-casein is
585 the precursor to its amyloid fibril formation. *Biochemical Journal*, 429, 251-260.
- 586 Farrell, Malin, E. L., Brown, C. J., & Qi, P. X. (2006). Casein micelle structure: What can be
587 learned from milk synthesis and structural biology? *Current Opinion in Colloid &*
588 *Interface Science*, 11, 135-147.
- 589 Farrell, H. M., Cooke, P. H., Wickham, E. D., Piotrowski, E. G., & Hoagland, P. D. (2003).
590 Environmental Influences on Bovine κ-Casein: Reduction and Conversion to Fibrillar
591 (Amyloid) Structures. *Journal of Protein Chemistry*, 22, 259-273.
- 592 Gasteiger, E., Hoogland, C., Gattiker, A., Duvaud, S., Wilkins, M. R., Appel, R. D., &
593 Bairoch, A. (2005). Protein Identification and Analysis Tools on the ExPASy Server.
594 In J. M. Walker (Ed.), *The Proteomics Protocols Handbook* (pp. 571-607). New
595 York, NY, USA: Humana.
- 596 Gebhardt, R., & Kulozik, U. (2014). Simulation of the shape and size of casein micelles in a
597 film state. *Food & Function*, 5, 780-785.
- 598 Ginger, M. R., & Grigor, M. R. (1999). Comparative aspects of milk caseins. *Comparative*
599 *Biochemistry and Physiology B-Biochemistry & Molecular Biology*, 124, 133-145.
- 600 Goldschmidt, L., Teng, P. K., Riek, R., & Eisenberg, D. (2010). Identifying the amyloome,
601 proteins capable of forming amyloid-like fibrils. *Proceedings of the National*
602 *Academy of Sciences*, 107, 3487-3492.
- 603 Habchi, J., & Longhi, S. (2012). Structural disorder within paramyxovirus nucleoproteins and
604 phosphoproteins. *Molecular Biosystems*, 8, 69-81.
- 605 Hambraeus, L., & Lonnerdal, B. (2003). Nutritional aspects of milk proteins. In P. F. Fox &
606 P. L. H. McSweeney (Eds.), *Advanced Dairy Chemistry* (Vol. 1 Part B, pp. 605-645).
607 New York, NY, USA: Kluwer Academic/Plenum.
- 608 Hodgkin, J. (1998). Seven types of pleiotropy. *International Journal of Developmental*
609 *Biology*, 42, 501-505.
- 610 Holt, C., & Sawyer, L. (1993). Caseins as Rheomorphic Proteins - Interpretation of Primary
611 and Secondary Structures of the Alpha-S1-Caseins, Beta-Caseins and Kappa-Caseins.
612 *Journal of the Chemical Society-Faraday Transactions*, 89, 2683-2692.
- 613 Holt, C., Wahlgren, N. M., & Drakenberg, T. (1996). Ability of a beta-casein phosphopeptide
614 to modulate the precipitation of calcium phosphate by forming amorphous dicalcium
615 phosphate nanoclusters. *Biochemical Journal*, 314, 1035-1039.
- 616 Holt, C., Timmins, P. A., Errington, N., & Leaver, J. (1998). A core-shell model of calcium
617 phosphate nanoclusters stabilized by beta-casein phosphopeptides, derived from
618 sedimentation equilibrium and small-angle X-ray and neutron-scattering
619 measurements. *European Journal of Biochemistry*, 252, 73-78.

- 620 Holt, C., de Kruif, C. G., Tuinier, R., & Timmins, P. A. (2003). Substructure of bovine casein
621 micelles by small-angle X-ray and neutron scattering. *Colloids and Surfaces a-
622 Physicochemical and Engineering Aspects*, 213, 275-284.
- 623 Holt, C., Sorensen, E. S., & Clegg, R. A. (2009). Role of calcium phosphate nanoclusters in
624 the control of calcification. *Fefs Journal*, 276, 2308-2323.
- 625 Holt, C., & Carver, J. A. (2012). Darwinian transformation of a 'scarcely nutritious fluid' into
626 milk. *Journal of Evolutionary Biology*, 25, 1253-1263.
- 627 Holt, C. (2013). Unfolded phosphoproteins enable soft and hard tissues to coexist in the same
628 organism with relative ease. *Current Opinion in Structural Biology*, 23, 420-425.
- 629 Holt, C., Carver, J. A., Ecroyd, H., & Thorn, D. C. (2013). Caseins and the casein micelle:
630 their biological functions, structures and behaviour in foods. *Journal of Dairy
631 Science*, 96, 6127–6146.
- 632 Holt, C., Lenton, S., Nylander, T., Sorensen, E. S., & Teixeira, S. C. M. (2014).
633 Mineralisation of soft and hard tissues and the stability of biofluids. *Journal of
634 Structural Biology*, 185, 383-396.
- 635 Huppertz, T., Kelly, A. L., & de Kruif, C. G. (2006). Disruption and reassociation of casein
636 micelles under high pressure. *Journal of Dairy Research*, 73, 294-298.
- 637 Ilawe, N. V., Raeber, A. E., Schweitzer-Stenner, R., Toal, S. E., & Wong, B. M. (2015).
638 Assessing backbone solvation effects in the conformational propensities of amino acid
639 residues in unfolded peptides. *Physical Chemistry Chemical Physics*, 17, 24917-
640 24924.
- 641 Ingham, B., Erlangga, G. D., Smialowska, A., Kirby, N. M., Wang, C., Matia-Merino, L.,
642 Haverkamp, R. G., & Carr, A. J. (2015). Solving the mystery of the internal structure
643 of casein micelles. *Soft Matter*, 11, 2723-2725.
- 644 Jackson, A. J., & McGillivray, D. J. (2011). Protein aggregate structure under high pressure.
645 *Chemical Communications*, 47, 487-489.
- 646 Jensen, M. R., Salmon, L., Nodet, G., & Blackledge, M. (2010). Defining Conformational
647 Ensembles of Intrinsically Disordered and Partially Folded Proteins Directly from
648 Chemical Shifts. *Journal of the American Chemical Society*, 132, 1270-1272.
- 649 Kalmar, L., Homola, D., Varga, G., & Tompa, P. (2012). Structural disorder in proteins
650 brings order to crystal growth in biominerilization. *Bone*, 51, 528-534.
- 651 Kawasaki, K., & Weiss, K. M. (2003). Mineralized tissue and vertebrate evolution: The
652 secretory calcium-binding phosphoprotein gene cluster. *Proceedings of the National
653 Academy of Sciences of the United States of America*, 100, 4060-4065.
- 654 Kawasaki, K., Lafont, A.-G., & Sire, J.-Y. (2011). The evolution of casein genes from tooth
655 genes before the origin of mammals. *Molecular Biology and Evolution*, 28, 2053-
656 2061.
- 657 Kay, B. K., Williamson, M. P., & Sudol, P. (2000). The importance of being proline: the
658 interaction of proline-rich motifs in signaling proteins with their cognate domains.
659 *Faseb Journal*, 14, 231-241.
- 660 Kyte, J., & Doolittle, R. F. (1982). A simple method for displaying the hydropathic character
661 of a protein. *Journal of Molecular Biology*, 157, 105-132.

- 662 Lefevre, C. M., Sharp, J. A., & Nicholas, K. R. (2009). Characterisation of monotreme
663 caseins reveals lineage-specific expansion of an ancestral casein locus in mammals.
664 *Reproduction Fertility and Development*, 21, 1015-1027.
- 665 Lefèvre, C. M., Sharp, J. A., & Nicholas, K. R. (2010). Evolution of lactation: ancient origin
666 and extreme adaptations of the lactation system. *Annual Review of Genomics and*
667 *Human Genetics*, 11, 219-238.
- 668 Lencki, R. W. (2007). Evidence for fibril-like structure in bovine casein micelles. *Journal of*
669 *Dairy Science*, 90, 75-89.
- 670 Leonil, J., Henry, G., Jouanneau, D., Delage, M.-M., Forge, V., & Putaux, J.-L. (2008).
671 Kinetics of Fibril Formation of Bovine κ -Casein Indicate a Conformational
672 Rearrangement as a Critical Step in the Process. *Journal of Molecular Biology*, 381,
673 1267-1280.
- 674 Liberles, D. A., Tisdell, M. D. M., & Grahnén, J. A. (2011). Binding constraints on the
675 evolution of enzymes and signalling proteins: the important role of negative
676 pleiotropy. *Proceedings of the Royal Society B-Biological Sciences*, 278, 1930-1935.
- 677 Lin, S. H. C., Dewan, R. K., Bloomfield, V. A., & Morr, C. V. (1971). Inelastic light-
678 scattering study of the size distribution of bovine milk casein micelles. *Biochemistry*,
679 10, 4788-4793.
- 680 Little, E. M., & Holt, C. (2004). An equilibrium thermodynamic model of the sequestration
681 of calcium phosphate by casein phosphopeptides. *European Biophysics Journal with*
682 *Biophysics Letters*, 33, 435-447.
- 683 Luck, G., Liao, H., Murray, N. J., Grimmer, H. R., Warminski, E. E., Williamson, M. P.,
684 Lilley, T. H., & Haslam, E. (1994). Polyphenols, astringency and proline-rich
685 proteins. *Phytochemistry*, 37, 357-371.
- 686 Marchin, S., Putaux, J.-L., Pignon, F., & Léonil, J. (2007). Effects of the environmental
687 factors on the casein micelle structure studied by cryo transmission electron
688 microscopy and small-angle x-ray scattering/ultrasmall-angle x-ray scattering. *J Chem*
689 *Phys*, 126, 045101.
- 690 Martin, P., Cebo, C., & Miranda, G. (2013). Interspecies comparison of milk proteins:
691 quantitative variability and molecular diversity. In P. L. H. McSweeney & P. F. Fox
692 (Eds.), *Advanced Dairy Chemistry* (4th ed., Vol. 1A. Proteins: Basic Aspects, pp. 386-
693 429). New York, NY, USA: Springer.
- 694 Mata, J. P., Udabage, P., & Gilbert, E. P. (2011). Structure of casein micelles in milk protein
695 concentrate powders via small angle X-ray scattering. *Soft Matter*, 7, 3837-3843.
- 696 McMahon, D. J., & McManus, W. R. (1998). Rethinking casein micelle structure using
697 electron microscopy. *Journal of Dairy Science*, 81, 2985-2993.
- 698 McMahon, D. J., & Oommen, B. S. (2012). Casein micelle structure, functions and
699 interactions. In P. F. Fox & P. L. H. McSweeney (Eds.), *Advanced Dairy Chemistry*
700 (4th ed., Vol. 1A Proteins: Basic Aspects, pp. 185-210). New York, NY, USA:
701 Springer.
- 702 Meral, D., Toal, S. E., Schweitzer-Stenner, R., & Urbanc, B. (2015). Water-Centered
703 Interpretation of Intrinsic pPII Propensities of Amino Acid Residues: In Vitro-Driven
704 Molecular Dynamics Study. *The Journal of Physical Chemistry B*, 119, 13237-13251.

- 705 Mercier, J. C. (1981). Phosphorylation of caseins, present evidence for an amino-acid triplet
706 code post-translationally recognized by specific kinases. *Biochimie*, 63, 1-17.
- 707 Mercier, J. C., & Villette, J. L. (1993). Structure and function of milk protein genes. *Journal*
708 *of Dairy Science*, 76, 3079-3098.
- 709 Murray, N. J., Williamson, M. P., Lilley, T. H., & Haslam, E. (1994). Study of the interaction
710 between salivary proline-rich proteins and a polyphenol by ¹H-NMR spectroscopy.
711 *European Journal of Biochemistry*, 219, 923-935.
- 712 Nelson, R., Sawaya, M. R., Balbirnie, M., Madsen, A. O., Riek, C., Grothe, R., &
713 Eisenberg, D. (2005). Structure of the cross-[beta] spine of amyloid-like fibrils.
714 *Nature*, 435, 773-778.
- 715 Niewold, T. A., Murphy, C. L., Hulskamp-Koch, C. A. M., Tooten, P. C. J., & Gruys, E.
716 (1999). Casein related amyloid, characterization of a new and unique amyloid protein
717 isolated from bovine corpora amylacea. *Amyloid-International Journal of*
718 *Experimental and Clinical Investigation*, 6, 244-249.
- 719 Oftedal, O. T. (2012). The evolution of milk secretion and its ancient origins. *Animal*, 6, 355-
720 368.
- 721 Ouanezar, M., Guyomarc'h, F., & Bouchoux, A. (2012). AFM Imaging of milk casein
722 micelles: Evidence for structural rearrangement upon acidification. *Langmuir*, 28,
723 4915-4919.
- 724 Pappu, R. V., Srinivasan, R., & Rose, G. D. (2000). The Flory isolated-pair hypothesis is not
725 valid for polypeptide chains: Implications for protein folding. *Proceedings of the*
726 *National Academy of Sciences of the United States of America*, 97, 12565-12570.
- 727 Pappu, R. V., Wang, X., Vitalis, A., & Crick, S. L. (2008). A polymer physics perspective on
728 driving forces and mechanisms for protein aggregation. *Archives of Biochemistry and*
729 *Biophysics*, 469, 132-141.
- 730 Pascal, C., Pate, F., Cheynier, V., & Delsuc, M.-A. (2009). Study of the Interactions Between
731 a Proline-Rich Protein and a Flavan-3-ol by NMR: Residual Structures in the Natively
732 Unfolded Protein Provides Anchorage Points for the Ligands. *Biopolymers*, 91, 745-
733 756.
- 734 Podgornaia, A. I., & Laub, M. T. (2015). Pervasive degeneracy and epistasis in a protein-
735 protein interface. *Science*, 347, 673-677.
- 736 Rijnkels, M., Kooiman, P. M., de Boer, H. A., & Pieper, F. R. (1997). Organization of the
737 bovine casein gene locus. *Mammalian Genome*, 8, 148-152.
- 738 Rodriguez, J. A., Ivanova, M. I., Sawaya, M. R., Cascio, D., Reyes, F. E., Shi, D., Sangwan,
739 S., Guenther, E. L., Johnson, L. M., Zhang, M., Jiang, L., Arbing, M. A., Nannenga,
740 B. L., Hattne, J., Whitelegge, J., Brewster, A. S., Messerschmidt, M., Boutet, S.,
741 Sauter, N. K., Gonen, T., & Eisenberg, D. S. (2015). Structure of the toxic core of
742 [agr]-synuclein from invisible crystals. *Nature*, 525, 486-490.
- 743 Sawaya, M. R., Sambashivan, S., Nelson, R., Ivanova, M. I., Sievers, S. A., Apostol, M. I.,
744 Thompson, M. J., Balbirnie, M., Wiltzius, J. J. W., McFarlane, H. T., Madsen, A. O.,
745 Riek, C., & Eisenberg, D. (2007). Atomic structures of amyloid cross-beta spines
746 reveal varied steric zippers. *Nature*, 447, 453-457.

- 747 Shukla, A., Narayanan, T., & Zanchi, D. (2009). Structure of casein micelles and their
748 complexation with tannins. *Soft Matter*, 5, 2884-2888.
- 749 Stewart, A. F., Bonsing, J., Beattie, C. W., Shah, F., Willis, I. M., & Mackinlay, A. G.
750 (1987). Complete nucleotide-sequences of bovine alpha-s2-casein and beta-casein
751 cDNAs - comparisons with related sequences in other species. *Molecular Biology and*
752 *Evolution*, 4, 231-241.
- 753 Tagliabracci, V. S., Engel, J. L., Wen, J., Wiley, S. E., Worby, C. A., Kinch, L. N., Xiao, J.,
754 Grishin, N. V., & Dixon, J. E. (2012). Secreted Kinase Phosphorylates Extracellular
755 Proteins That Regulate Biomineralization. *Science*, 336, 1150-1153.
- 756 Theillet, F.-X., Kalmar, L., Tompa, P., Han, K.-H., Selenko, P., Dunker, A. K., Daughdrill,
757 G. W., & Uversky, V. N. (2013). The alphabet of intrinsic disorder. *Intrinsically*
758 *Disordered Proteins*, 1, e24360.
- 759 Thompson, M. J., Sievers, S. A., Karanicolas, J., Ivanova, M. I., Baker, D., & Eisenberg, D.
760 (2006). The 3D profile method for identifying fibril-forming segments of proteins.
761 *Proceedings of the National Academy of Sciences of the United States of America*,
762 103, 4074-4078.
- 763 Thorn, D. C., Meehan, S., Sunde, M., Rekas, A., Gras, S. L., MacPhee, C. E., Dobson, C. M.,
764 Wilson, M. R., & Carver, J. A. (2005). Amyloid fibril formation by bovine milk
765 kappa-casein and its inhibition by the molecular chaperones alpha(s-) and beta-casein.
766 *Biochemistry*, 44, 17027-17036.
- 767 Thorn, D. C., Ecroyd, H., Sunde, M., Poon, S., & Carver, J. A. (2008). Amyloid fibril
768 formation by bovine milk alpha(s2)-casein occurs under physiological conditions yet
769 is prevented by its natural counterpart, alpha(s1)-casein. *Biochemistry*, 47, 3926-3936.
- 770 Thorn, D. C., Ecroyd, H., Carver, J. A., & Holt, C. (2015). Casein structures in the context of
771 unfolded proteins. *International Dairy Journal*, 46, 2-11.
- 772 Tokuriki, N., Stricher, F., Serrano, L., & Tawfik, D. S. (2008). How Protein Stability and
773 New Functions Trade Off. *Plos Computational Biology*, 4, Article Number: e1000002
- 774 Tran, H. T., Mao, A., & Pappu, R. V. (2008). Role of Backbone-Solvent Interactions in
775 Determining Conformational Equilibria of Intrinsically Disordered Proteins. *Journal*
776 *of the American Chemical Society*, 130, 7380-7392.
- 777 Trejo, R., Dokland, T., Jurat-Fuentes, J., & Harte, F. (2011). Cryo-transmission electron
778 tomography of native casein micelles from bovine milk. *Journal of Dairy Science*, 94,
779 5770-5775.
- 780 Treweek, T. M., Thorn, D. C., Price, W. E., & Carver, J. A. (2011). The chaperone action of
781 bovine milk alpha(S1)- and alpha(S2)-caseins and their associated form alpha(S)-
782 casein. *Archives of Biochemistry and Biophysics*, 510, 42-52.
- 783 Umeda, T., & Aoki, T. (2002). Relation between micelle size and micellar calcium
784 phosphate. *Milchwissenschaft-Milk Science International*, 57, 131-133.
- 785 Uversky, V. N., Gillespie, J. R., & Fink, A. L. (2000). Why are "natively unfolded" proteins
786 unstructured under physiologic conditions? *Proteins-Structure Function and Genetics*,
787 41, 415-427.

- 788 Uversky, V. N., Oldfield, C. J., & Dunker, A. K. (2005). Showing your ID: intrinsic disorder
789 as an ID for recognition, regulation and cell signaling. *Journal of Molecular*
790 *Recognition*, 18, 343-384.
- 791 Wang, Z., Liao, B.-Y., & Zhang, J. (2010). Genomic patterns of pleiotropy and the evolution
792 of complexity. *Proceedings of the National Academy of Sciences of the United States*
793 *of America*, 107, 18034-18039.
- 794 Weinreich, D. M., Delaney, N. F., DePristo, M. A., & Hartl, D. L. (2006). Darwinian
795 evolution can follow only very few mutational paths to fitter proteins. *Science*, 312,
796 111-114.
- 797 Westerheide, S. D., Raynes, R., Powell, C., Xue, B., & Uversky, V. N. (2012). HSF
798 Transcription Factor Family, Heat Shock Response, and Protein Intrinsic Disorder.
799 *Current Protein & Peptide Science*, 13, 86-103.
- 800 Williamson, M. P. (1994). The structure and function of proline-rich regions in proteins.
801 *Biochemical Journal*, 297, 249-260.
- 802 Wiltzius, J. J. W., Sievers, S. A., Sawaya, M. R., Cascio, D., Popov, D., Riek, C., &
803 Eisenberg, D. (2008). Atomic structure of the cross- β spine of islet amyloid
804 polypeptide (amylin). *Protein Science*, 17, 1467-1474.
- 805 Zarrinpar, A., Park, S. H., & Lim, W. A. (2003). Optimization of specificity in a cellular
806 protein interaction network by negative selection. *Nature*, 426, 676-680.
- 807 Zeng, Y., & Gu, X. (2010). Genome factor and gene pleiotropy hypotheses in protein
808 evolution. *Biology Direct*, 5, Article Number: 37
- 809

Supplementary Data

Casein micelle structures, functions and diversity C. Holt (2016) Int Dairy J.

Data S1. Accession codes of sequence data

P02662	CASA1_BOVIN	Alpha-S1-casein	CSN1S1	Bos taurus (Bovine)
O62823	CASA1_BUBBU	Alpha-S1-casein	CSN1S1	Bubalus bubalis (Domestic water buffalo)
P18626	CASA1_CAPHI	Alpha-S1-casein	CSN1S1	Capra hircus (Goat)
P04653	CASA1_SHEEP	Alpha-S1-casein	CSN1S1	Ovis aries (Sheep)
CX987842	CSN1S1_DOG	Alpha-S1-casein	CSN1S1	Canis lupus familiaris
P39035	CASA1_PIG	Alpha-S1-casein	CSN1S1	Sus scrofa (Pig)
H2QPK8	H2QPK8_PANTR	Alpha-S1-casein	Uncharacterized protein	Pan troglodytes (Chimpanzee)
O97943	CASA1_CAMDR	Alpha-S1-casein	CSN1S1	Camelus dromedarius (Dromedary) (Arabian camel)
C3W972	C3W972_EQUAS	Alpha s1 casein	csn1S1	Equus asinus africanus (donkey)
P47710	CASA1_HUMAN	Alpha-S1-casein	CSN1S1 CASA CSN1	Homo sapiens (Human)
Q9XSE3	Q9XSE3_TRIVU	Alpha-casein		Trichosurus vulpecula (Brush-tailed possum)
D0QJ96	D0QJ96_ORNAN	Alpha casein	CSN1	Ornithorhynchus anatinus (Duckbill platypus)
D0QJA2	D0QJA2_9MAMM	Alpha casein	CSN1	Tachyglossus aculeatus (Australian echidna)
P09115	CASA1_RABIT	Alpha-S1-casein	CSN1S1	Oryctolagus cuniculus (Rabbit)
P19228	CASA1_MOUSE	Alpha-S1-casein	Csn1s1 Csn1 Csna	Mus musculus (Mouse)
P04656	CASA1_CAVPO	Alpha-S1-casein	CSN1S1	Cavia porcellus (Guinea pig)
P02663	CASA2_BOVIN	Alpha-S2-casein	CSN1S2	Bos taurus (Bovine)
B6VPY2	B6VPY2_BUBBU	Alpha s2 casein	csn1s2	Bubalus bubalis (Domestic water buffalo)
P33049	CASA2_CAPHI	Alpha-S2-casein	CSN1S2	Capra hircus (Goat)
P04654	CASA2_SHEEP	Alpha-S2-casein	CSN1S2	Ovis aries (Sheep)
P39036	CASA2_PIG	Alpha-S2-casein	CSN1S2	Sus scrofa (Pig)
O97944	CASA2_CAMDR	Alpha-S2-casein	CSN1S2	Camelus dromedarius (Dromedary) (Arabian camel)
B7VGF9	CASA2_EQUAS	Alpha-S2-casein	CSN1S2	Equus asinus (Donkey)
P04655	CASA2_CAVPO	Alpha-S2-casein	CSN1S2	Cavia porcellus (Guinea pig)
P50419	CASA2_RABIT	Alpha-S2-casein	CSN1S2	Oryctolagus cuniculus (Rabbit)
Q02862	CS2LA_MOUSE	Alpha-S2-casein-like A	Csn1s2a Csng	Mus musculus (Mouse)
P02667	CS2LA_RAT	Alpha-S2-casein-like A	Csn1s2a Csng	Rattus norvegicus (Rat)
P50418	CS2LA_RABIT	Alpha-S2-casein-like A	CSN1S2A	Oryctolagus cuniculus (Rabbit)
P02664	CS2LB_MOUSE	Alpha-S2-casein-like B	Csn1s2b Csnd Csne	Mus musculus (Mouse)

Q8CGR3	CS2LB_RAT	Alpha-S2-casein-like B	Csn1s2b Csnd	Rattus norvegicus (Rat)
D0QJ98	D0QJ98_ORNAN	Beta-like casein 2		Ornithorhynchus anatinus (Duckbill platypus)
D0QJA6	D0QJA6_9MAMM	Beta-like casein 2 variant 3		Tachyglossus aculeatus (Australian echidna)
P02666	CASB_BOVIN	Beta-casein	CSN2	Bos taurus (Bovine)
Q9TSI0	CASB_BUBBU	Beta-casein	CSN2	Bubalus bubalis (Domestic water buffalo)
P33048	CASB_CAPHI	Beta-casein	CSN2	Capra hircus (Goat)
P11839	CASB_SHEEP	Beta-casein	CSN2	Ovis aries (Sheep)
P39037	CASB_PIG	Beta-casein	CSN2	Sus scrofa (Pig)
Q9TVD0	CASB_CAMDR	Beta-casein	CSN2	Camelus dromedarius (Dromedary) (Arabian camel)
Q9GKK3	CASB_HORSE	Beta-casein	CSN2	Equus caballus (Horse)
Q9N2G8	Q9N2G8_CANFA	Beta-casein		Canis familiaris (Dog) (Canis lupus familiaris)
P05814	CASB_HUMAN	Beta-casein	CSN2 CASB	Homo sapiens (Human)
P09116	CASB_RABIT	Beta-casein	CSN2	Oryctolagus cuniculus (Rabbit)
P02665	CASB_RAT	Beta-casein	Csn2 Csnb	Rattus norvegicus (Rat)
P10598	CASB_MOUSE	Beta-casein	Csn2 Csnb	Mus musculus (Mouse)
G3U197	G3U197_LOXAF	Uncharacterized protein	CSN2	Loxodonta africana (African elephant)
Q9XSE4	CASB_TRIVU	Beta-casein	CSN2 BCAS	Trichosurus vulpecula (Brush-tailed possum)
D0QJ95	D0QJ95_ORNAN	Beta casein	CSN2	Ornithorhynchus anatinus (Duckbill platypus)
D0QJ99	D0QJ99_9MAMM	Beta casein	CSN2	Tachyglossus aculeatus (Australian echidna)
D0QJA4	D0QJA4_9MAMM	Beta-like casein 2 variant 1		Tachyglossus aculeatus (Australian echidna)
P02668	CASK_BOVIN	Kappa-casein	CSN3 CSN10 CSNK	Bos taurus (Bovine)
P11840	CASK_BUBBU	Kappa-casein	CSN3 CSN10 CSNK	Bubalus bubalis (Domestic water buffalo)
P02669	CASK_SHEEP	Kappa-casein	CSN3 CSN10 CSNK	Ovis aries (Sheep)
P02670	CASK_CAPHI	Kappa-casein	CSN3 CSN10 CSNK	Capra hircus (Goat)
P79139	CASK_CAMDR	Kappa-casein	CSN3 CSN10 CSNK	Camelus dromedarius (Dromedary) (Arabian camel)
P11841	CASK_PIG	Kappa-casein	CSN3 CSN10 CSNK	Sus scrofa (Pig)
P33618	CASK_RABIT	Kappa-casein	CSN3 CSN10 CSNK	Oryctolagus cuniculus (Rabbit)
P19442	CASK_CAVPO	Kappa-casein	CSN3 CSN10 CSNK	Cavia porcellus (Guinea pig)
P04468	CASK_RAT	Kappa-casein	Csn3 Csn10 Csnk	Rattus norvegicus (Rat)
P06796	CASK_MOUSE	Kappa-casein	Csn3 Csn10 Csnk	Mus musculus (Mouse)
P07498	CASK_HUMAN	Kappa-casein	CSN3 CASK CSN10	Homo sapiens (Human)
P82187	CASK_HORSE	Kappa-casein	CSN3 CSN10	Equus caballus (Horse)
E2QXF8	E2QXF8_CANFA	Uncharacterized protein	CSN3	Canis familiaris (Dog) (Canis lupus familiaris)
G3UDT9	G3UDT9_LOXAF	Uncharacterized protein	CSN3	Loxodonta africana (African elephant)

D0QJA9	D0QJA9_ORNAN	Kappa casein	CSN3	Ornithorhynchus anatinus (Duckbill platypus)
D0QJA7	D0QJA7_9MAMM	Kappa casein	CSN3	Tachyglossus aculeatus (Australian echidna)
Q9XSD6	Q9XSD6_TRIVU	Kappa casein	CASK	Trichosurus vulpecula (Brush-tailed possum)
F7E1V6			CSN3	Monodelphis domestica (Gray-tailed opossum)

Multiple Sequence alignments

Multiple amino acid sequence alignments α_{S1} -casein (CSN1S1; Data S2), α_{S2} -casein (CSN1S2; Data S3), α_{S2} -casein B-type (CSN1S2B; Data S4), β -casein (CSN2; Data S5) and κ -casein (CSN3; Data S6) were based on the previous work of Kawasaki *et al.*, (Kawasaki, Lafont, & Sire, 2011) but with an increased number of sequences, alignments of sequences into 5 rather than 4 groups and a refined method of alignment. Monotreme CSN2B amino acid sequences coded by exons 2-6 were aligned with CSN2, whereas those coded by exon 7 and the following exons were aligned with CSN1S2B. Polar tract sequences are coloured green. Canonical sites of phosphorylation by the Golgi kinase are coloured red. Identities are denoted by # in the last row and indels by -. Boxed sequences are encoded by candidate orthologous exons.

Data S2. Alignment of amino acid sequences of CSN1S1 genes

cow	MKLLILTCVALVALARP	-----	-----	KHPIKHQGL--PQ-	-----	-----	EVLNE-NLLRFFVA-----	PFPEVFGK	EKVNELSK
buffalo	MKLLILTCVALVALARP	-----	-----	KQPIKHQGL--PQ-	-----	-----	GVLNE-NLLRFFVA-----	PFPEVFGK	EKVNELST
goat	MKLLILTCVALVALARP	-----	-----	KHPINHRGL--SP-	-----	-----	EVPNE-NLLRFVVA-----	PFPEVFRK	ENINELSK
sheep	MKLLILTCVALVALARP	-----	-----	KHPIKHQGL--SS-	-----	-----	EVLNE-NLLRFVVA-----	PFPEVFRK	ENINELSK
pig	MKLLIFICLAAVALARP	-----	-----	KPPLRHQEHH--LQ-	-----	NE----PDSRE	ELFKERKFLRFPEV-----	PILLSQFRQ	EIINELNR
camel	MKLLILTCVALVALARP	-----	-----	KYPLRYPEV--FQ-	-----	NE----PDSIE	EVLNK-RKILELAV-----	VSPIQFRQ	ENIDEL-K
donkey	MKLLILTCVALVALARP	-----	-----	KLPHRHPET--IQ-	-----	NE----QDSRE	KVLKE-RKFPSFAL-----	HT--SRE	EYINELNR
dog	MKFLLITCLVALVALARP	-----	-----	KLPLRHPET--TQ-	-----	NE----LDSRE	EVLKERQFLRF-AL-----	PTPRELRE	-----
human	MRLLLITCLVALVALARP	-----	-----	KLPLRYPER--LQ-	-----	NP----SESSE	-----	PIPELSRE	EYMNMGMR
chimpanzee	MRLLLITCLVALVALARP	-----	-----	KLPLRYPER--LQ-	-----	NP----SESSE	-----	PIPELSRQ	EYMNMGMR
guinea pig	MKLLILTCVALVASAVAMP	-----	-----	KFPFRTEL--FQ-	-----	TQRGGSSSSSSSE	ERLKE-ENIFKFDQ-----	QKELQ-RK	-----
rabbit	MKLLILTCVALTALARH	-----	-----	KFHGLHLKL--TQ-	-----	EQP----ESSEQ	EILKE-RKLLRFVQ-----	TVPLELRE	EYVNELNR
mouse	MKLLILTCVALAAAFAMP	-----	-----	RLHSRNASV--SQ-	-----	TQQQ----HSSSE	EIFKQ-PKYLNLNQ-----	--DLRQ	EFVNNMNR
possum	MKLLIFSCLMALALARE	-----	DVLHLSID	R-HIKHREVENRS-	NEDLIPLNE	-----VSSSE	ESLHQLNRDRRSPEKYELNKYRE	-----	-----
platypus	MKVLLACLVAVAVAMP	ESPSSSSSE	EAPRLLTK	KRILRNQEYYLPHL	-----	EE--SRSSSSSE	-----	ESTRPTLK	ESTDRDLKR
echidna	MKVLLACLVAVVAMP	ESPSSSSSE	EASKILTK	KRVQRDQEYYLPHQ	-----	EE--SVSSSSSE	-----	-----	-----
	# ## # # #								

cow	-----	DIGSESTE	DQAMEDIK	QMEAESISSSE	E---IVPNSVE	-----	QKHIQK-EDVPSERYLGYL	EQLRLKKYKVPQL	EIVPNSAE	-----
buffalo	-----	DIGSESTE	DQAMEDIK	QMEAESISSSE	E---IVPISVE	-----	QKHIQK-EDVPSERYLGYL	EQLRLKKYKVPQL	EIVPNLAE	-----
goat	-----	DIGSESTE	DQAMEDAK	QMKAGSSSSSE	E---IVPNSAE	-----	QKYIQK-EDVPSERYLGYL	EQLRLKKYKVPQL	EIVPKSAE	-----
sheep	-----	DIGSESIE	DQAMEDAK	QMKAGSSSSSE	E---IVPNSAE	-----	QKYIQK-EDVPSERYLGYL	EQLRLKKYKVPQL	EIVPKSAE	-----
pig	-----	-----	NHGMGHE	Q-RGSSSSSE	E---VVGNSAE	-----	QKHVQKEEDVPSQSYL---	GHLQGLNKYKLRQL	-----	EAIHDQ
camel	-----	DTRNEPTE	DHIMEDTE	RKESG-SSSSE	E---VVSSTTE	-----	QKDILK-EDMPSQRYL--	EELHRLNKYKLLQL	-----	EAIRDQ
donkey	QRELLKEKQKDEHK	-----	EYLIEDPE	QQESSSTSSSE	E---VVPINTE	-----	QKRIPR-EDMLYQHTL--	EALRRLSKYKQNLQL	-----	QAIYAQ
dog	-RELLREKQNEGIK	-----	-----	QRQSSSTSSSE	E---VVPNNTE	-----	QRQIPR-EDILYQRYL--	EQLRRLSQHNQQLQ-	-----	GTIHDO
human	QRNILREKQTDEIK	DTRNESTQ	NCVVAEPE	KMESSSISSSE	E---MSLSKCA	-----	-----	EQFCRLNEYQNQQLQ	-----	QAHAHQ
chimpanzee	QRNILREKQTDEIK	DTRNESTQ	NCVMAEPE	KMESSSISSSE	E---ISLSKCA	-----	-----	EQLSRLIKYHQLM	-----	QAVHAQ
guinea pig	QS--EKIK	EIISESTE	-----	QREASSSISSSE	E---VVPKNT	-----	QKHIPQ-EDALYQQAL--	-----	-----	EVVHAQ
rabbit	QRELLREKENEEIK	GTRNEVTE	EHVLADRE	-TEASISSSE	E---IVPSSTK	-----	QKYVPR-EDLAYQPYV--	-----	-----	-----
mouse	QRALLTE-QNDEIK	VTMDAASE	EQAMASAQ	E-DSSSISSSE	ESEEAIPNITE	-----	QKNIAN-EDMLNQCTL--	EQLQROFKYNQLLQ	-----	KASLAK
possum	-----	-----	-----	DLKT--SSSE	ES--VAP-STE	ESVRQ	VEYNFN-EQEDASASRE--	RKIEDVSEQYRQYL	-----	-----
platypus	-RLLLKEKPILEHIL	-----	-----	KAPE--SSSSE	ES---DSAEE	-----	KRLLR--EREFYQQQL--	-----	-----	-----
echidna	-R-LLKDKPIFRLL	-----	-----	KATE--SSSSE	ES---DSAIE	-----	KRILR--ERQYYQQKL--	-----	-----	-----
				####	#					

cow	-ERLHSMKEGIHAQQ	KEPMIGVNQ	ELAYFYPE	-----LFRQFYQLDAYPSGAWWYVPLGTQYTDAPSFSIDPNPIGSENSEKTT-MPLW-----	
buffalo	-EQLHSMKEGIHAQQ	KEPMIGVNQ	ELAYFYPQ	-----LFRQFYQLDAYPSGAWWYVPLGTQYDAPSFSIDPNPIGSENSGKTT-MPLW-----	
goat	-EQLHSMKEGNPAHQ	KQPMIAVNQ	ELAYFYPQ	-----LFRQFYQLDAYPSGAWWYLPPLGTQYTDAPSFSIDPNPIGSENSGKTT-MPLW-----	
sheep	-EQLHSMKEGNPAHQ	KQPMIAVNQ	ELAYFYPQ	-----LFRQFYQLDAYPSGAWWYLPPLGTQYTDAPSFSIDPNPIGSENSGKIT-MPLW-----	
pig	-E-LHRTNEDKHTQQ	GEPMKGVNQ	EQAYFYFE	-----PLHQFYQLDAYPYATWYYPP---QYIAHPLFTNIPQPTAPEKGKTEIMPOW-----	
camel	-EQYLRYINEDNHQPQL	GEPVKVVVTQ	EQAYFHLE	-----PFPQFFQLGASPYVAWWYPPQVMQYIAHPSSYDTPPEGIASEDGKTDVMPQWW-----	
donkey	-EQLLRMKNS--QR	K-PMRVRVNQ	-----	-----PFQPSYQLDVVPYAAWFHBAQMHQHVAYSFHDTKLIASENSEKTDIPEW-----	
dog	-QQLLRRVNNENLLQL	-----	-----	-----PFQFYQLDAYPYAWFPAQIMQYIAYPPLSDITKPIASENIEADNVVPQW-----	
human	-EQIRRNMENSHVQV	-----	-----	-----PFQQLNQLAAYPYAVWYY-PQIMQYVPPFSIDSNPTAHENYEKNVMLQW-----	
chimpanzee	-EQIRRNMENSHVQV	-----	-----	-----PFQQLNQLAAYPCAVWYY-PQIMQYVPPFSIDSNPTAHENYEKNVMLQW-----	
guinea pig	-EQFHRINEHNQAQV	KEPMR VFNQ	-----	-----LDAYPFAAWYYGPE-VQYMSFLPFSIOPQIFPEDAQNTEVMPPEWVM-----	
rabbit	QQQLLRRMKERYQIQE	REPMRVVNQ	E LAQLYLQ	-----PFEQPYQLDAYLAPWYYTPEVMQYVLSPLFYDLVTPSAFE SAEKTDVPIPEWLKN-----	
mouse	-EQPYRMNAYSQVQM	RHPMSVVDQ	ALA QFSVQ	-----PFPQIFQYDAFPL---WAYFPQDMQYLT PKA VLN TFKP IVSKD TEKTNV---W-----	
possum	-----	-----	EPL -YYAT	-----EP-DFYYTIVPISMPRF FPYPAEAPVFSTRKAPVPSINRATEAVYTY SEEK-----	KN
platypus	-----	-----	-----	D---EYYRQFEP-DFYPRAYPKK--EVMPYPLEYFIPQAAVYSI PQLVYRV PQEVTFPSPSLRFYAFPQPTLPVE	RK
echidna	-----	-----	-----	DELKEYFRQFEP-YFYPVAYQKK--EVMPYQLEYFVPQPEVY S I P QPVYRV PQEVTFPSPSLHFYAFPQSTLPIE	RK

Data S3. Alignment of amino acid sequences from CSN1S2 genes

cow	MKFFIFTCCLLAVALAKH	TMEHVSSSE	ESI--ISQE----	TYKQEKNMIAHPSK	ENLCSTFCK	EVVRNA--NEE	-----	EYSIGSSSE	ESAEVATE
buffalo	MKFFIFTCCLLAVALAKH	TMEHVSSSE	ESI--ISQE----	TYKQEKNMIAHPSK	ENLCSTFCK	EVIRNA--NEE	-----	EYSIGSSSE	ESAEVATE
goat	MKFFIFTCCLLAVALAKH	KMEEHVSSSE	EPI--NIFQE----	TYKQEKNMIAHPRK	EKLCTTSCE	EVVRNA--NEE	-----	EYSIRSSSE	ESAEVAPE
sheep	MKFFIFTCCLLAVALAKH	KMEEHVSSSE	EPI--NISQE----	TYKQEKNMIAHPRK	EKLCTTSCE	EVVRNA--DEE	-----	EYSIRSSSE	ESAEVAPE
pig	MKFFIFTCCLLAVALAFAKH	EMEHVSSSE	ESI--NISQE----	KYKQEKNVINHPSK	EDICATSC	EAVRNI--KEV	-----	GYASSSSSE	ESVDIPAE
camel	MKFFIFTCCLLAVALAKH	EMDGQSSE	ESI--NVSQQ----	KFKQVKKVAIHPSK	EDICSTFCE	EAVRNI--KEV	-----	-----	ESAEPVTE
donkey	MKFFIFTCCLLAVALAKH	NMEHRSSE	DSV--NISQE----	KFKQEKYVVIPTSK	ESICSTSC	EATRNI--NEM	ESAKFPTE	VYSSSSSE	ESAKFPTE
guinea pig	MKLFIFTCLLAVALAKH	KSEQQSSE	ESV--SISQE----	KFK--DKNMDTISSE	ETICASLCK	EATKNT--PKM	-----	AFFSRSSSE	EFADIHRE
rabbit a	MRFFVFTCLLAVALAKH	GIEQRSSAE	EIV--SFYQE----	KYKQDSNAAIYPTN	-----	-----	-----	--SVSSSE	ESVEVQTE
mouse a	MKFFIFACLVVVALAKH	EIKDKSSSE	ESSASIYPG----	KSKLDNSVFFQTTK	-----	-----	-----	DSASSSSSE	ESSEEVSE
rat a	MKFFIFTCCLVAAALAKH	AVDKDPSE	ESA--SVYLG----	KYKQGNGVFFQTPQ	-----	-----	-----	DSASSSSSE	ESSEEISE
	# # # ##	##	#					# #	

cow	EVKITVDDKHYQKAL	NEINQFYQK--FPQYLQYLYQG-PIVLNPWDQVKRNAVPI-TPTL	NR--EQLSTSE	-----	ENSKKTVDM	-----
buffalo	EVKITVDDKHYQKAL	NEINQFYQK--FPQYLQYLYQG-PIVLNPWDQVKRNAVPI-TPTL	NR--EQLSTSE	-----	ENSKKTVDM	-----
goat	EIKITVDDKHYQKAL	NEINQFYQK--FPQYLQYLYQG-PIVLNPWDQVKRNAVPI-TPTL	NR--EQLSTSE	-----	ENSKKTIDM	-----
sheep	EVKITVDDKHYQKAL	NEINQFYQK--FPQYLQYLYQG-PIVLNPWDQVKRNAVPI-TPTV	NR--EQLSTSE	-----	ENSKKTIDM	-----
pig	NVKVTVEDKHYLKQL	NEINQFYQK--FPQYLQYLYQG-PIVLNPWDQVKRNAVPI-TPTV	IQSGEELSTSE	-EPVSSSQE	E-NTKTVDM	-----
camel	-----	EKISQFYQK--FPQYLQALYQA-QIVMNPWDQTKTSAYPF-IPTV	NT--EQLSISE	-----	E-STEVPTE	-----
donkey	REEKEVEEKHHHLKQL	NKISQFYQKWKFLQYLQALHOG-QIVMNPWDQGKTRAYPF-IPTV	NT--EQLFTSE	-----	EIPKKTVDM	-----
guinea pig	-----	NKINQFYEKLNFQYLQALRQP-RIVLTPWDQTKTGAASP-IPIV	GK--EQISTIE	-----	DILKTTAV	ESSSSSSTE
rabbit a	KDEQIEEENVYLKQL	NKKDQLYQKWMVPQYNPDFYQR-PVVMSPWNQIYTRPYPIVLPTL	-----	PETTRIPLE	EIVKKIVEM	-----
mouse a	KIVQSEEQKVNLNQQ	KRIKQIFQKFYIPQYVE-VYQQ-QIVMNPWVKVTTTYPV--PI-	-----	-ESISTSVE	EILKKIIDM	-----
rat a	KIEQSEEQKVNLNQQ	KKFKQFSQESSFSQCCTPLHQQQQSSSVNQWPQP--NAIHN-TPTQ	-----	-ESTSTSVE	EILKKIDI	-----
		KKSQFSQDSSFPQICT-PYQQ-QSSVNQWPQP--NAIYD-VPSQ	-----			
	# # # # #	# # # # #	#			

cow	E STEVFTK	KTKLTEEKNRLNFL	KKISQRYQKPALPQYLKTVYQHQKAMKPWIOPKTKV----IPYV	RYL		
buffalo	E STEVITK	KTKLTEEDKNRLNFL	KKISQHYQKFTWPQYLKTVYQHQKAMKPWTQPKTNV----IPYV	RYL		
goat	ESTKVFHK	KTKLTEEKNRLNFL	KIIISQYYQKFAWPQYLKTVDQHQKAMKRWTQPKTN-----IPYV	RYL		
sheep	E STEVFTK	KTKLTEEKNRLNFL	KKISQYYQKFAWPQYLKTVDQHQKAMKPWTQPKTN-----IPYV	RYL		
pig	E SMEEFTK	KTELTEEKNRIKFL	NKIKQYYQKFTWPQYLKTVHQQKQAMKPWNHIKTNSYQI--IPNL	RYF		
camel	E STEVFTK	KTELTEEKNRLNFL	NKIQYQQTFWLPEYIPLTVYQYQKTMTPWNHIK-----	RYF		
donkey	E STEVVTE	KTELTEEKNYLKLL	NKINQYYEKFTLPQYFKIVHQHQTMDPQSHSKTNSYQI--IPVL	RYF		
guinea pig	KSTDVFHK	KTKMDEVQKLIQSSL	NIIHEYSQKAFWSQTLLEDVDQYLFKVFMPWNHYNTNADQVD-ASQE	RQA		
rabbit a	-----	-----	IKFNQ-LHQFVIPQYVQALQQ-RIAMNPWHHVTFRS----FPV-	LNF		
mouse a	-----	-----	IKYIQ-YQQVTIPQLPQALHP-QIPVSYWPSKDYTFPN--AHYT	RFY		
rat a	-----	-----	VKYFQ-YQQLTNPHFPQAVHP-QIPVSSWAPSKDYTFPT--ARYM	--A		

Data S4. Alignment of CSN1S2B-like sequences

rabbit b	MKFFIFTCLLAVALAVALAKP	KIEQ- SSSE	ETI-AVSQEVS P NL	-----	ENICSTACE	EPIKNI--NEV	EYVEVPTE
mouse b	MKFII L TCLLAVALAVALAQ	RMEQY ISSE	ESM-DNSQE-----	NFKQNMDVAFFP SQ	-----	ETVENIYIPQM	E SVEAPMK
rat b	MKFII L TCLLAVALAVALAQ	-----	ESK-DNSQE-----	DFKQTVDVVIFPGQ	-----	ETVKNIPIPQM	E SVEAPIK

platypus 2b	IKDQE FY QKVNL L QYLQALYQY-PTVMDPWTRAETKAIPF-IRTM	-----	KYQQR-LRLFKPTYL	VPVNKFVE-	RHPFRNILFPEELPEAYQPIE
echidna 2b	-----	-VSDII SQ	IYQQG-LRPFKPTHL	-----	RRPLKYIFF SE PPKVVQPIQ
rabbit b	NKCYQS I QTFKPPQALKGLYQY-HMAKNP WGY TVNRAFPS-TRTL	-----	QYKQE K -DATKHTSQ	-----	-----
mouse b	-----	-----	QYNQKMM DMS V S A R E	-----	-----
rat b	-----	-----	QYNQKTMDLSMRARE	-----	-----

platypus 2b	KED SSSSSE	ETVQVPVE	K-HLLRLRK-LHVPQ	-KLRP--LRFYPNHQVPFQRHPLPYAG----TQVHQPVEVPFPLP	VQY
echidna 2b	NED SSSSSE	EPVEVPAE	QNHLRLKK-LQVLQ	-NLQP--LRLLPNYQVPLQRHPLPVRLPNVFQAPHV EL PFPLP	QVV
rabbit b	-----	-----	KTELTEEEKAFLKYL	DEM K QYYQKFVF P QYLKNAHHFQ K TMNPWNHVKTIIYQS--VPTL	RYL
mouse b	-----	-----	KTVMTEESKNI Q DYM	NKMKR-YSKITWPQFVKLLHQYQKTMT P WSYY P ST-----PSQ	--V
rat b	-----	-----	KIVMSEIKKNI Q DYV	TKMKQ-YSKITWP R FVKSLQQYQKT M NPW S C P YT-----LLQ	--V

Data S5. Alignment of amino acid sequences of CSN2 genes

cow	MKVLLILACLVALALARE---	LEELNVPGE	IVE SL -----SSSE	-----	-----	-----	-----	ESITRINK	-KIEKFQ S EEQQQTE
buffalo	MKVLLILACLVALALARE---	LEELNVPGE	IVE SL -----SSSE	-----	-----	-----	-----	ESITHINK	-KIEKFQ S EEQQQME
goat	MKVLLILACLVALAIARE---	QEELNVVGE	TVE SL -----SSSE	-----	-----	-----	-----	ESITHINK	-KIEKFQ S EEQQQTE
sheep	MKVLLILACLVALALARE---	QEELNVVGE	TVE SL -----SSSE	-----	-----	-----	-----	ESITHINK	-KIEKFQ S EEQQQTE
pig	MKLLLILACFVALALARA---	KEELNASGE	TVE SL -----SSSE	-----	-----	-----	-----	ESITHISK	EKIEKLKREEQQOTE
camel	MKVLLILACRVALALARE---	KEEFKTAGE	ALESI-----SSSE	-----	-----	-----	-----	ESITHINK	QKIEKFKIEEQQOTE
horse	MKLLILACLVALALARE---	KEELNVSSE	TVESLSSNEPDSSSE	-----	-----	-----	-----	-----	EKLQKFKHEQQORE
dog	MKVFLILACLVALALARE---	KEELTLSNE	TVE SL -----SSSE	-----	-----	-----	-----	ESITHINK	OKLENFKHHEQQORE
human	MKVLLILACLVALALARE---	-----	TIESL-----SSSE	-----	-----	-----	-----	ESITEY-K	OKVEVKVKHEDQQGEG
rabbit	MKVLLILACLVALALARE---	KEQLSVPTE	AVGSV-----SSSE	-----	-----	-----	-----	E-IITHINK	QKLETIKHVQEQLRE
rat	MKVFLILACLVALALARE---	KDAFTVSSE	T-GSI-----SSE	-----	-----	-----	-----	ESVEHINE	-KLQKVKLMLGQVQSE
mouse	MKVFLILACLVALALARE---	-TTFTVSSE	T-DSI-----SSE	-----	-----	-----	-----	ESVEHINE	OKLQKVNLMGQLOAE
elephant	MKVFLILACLVAFALGRE---	KEEIIIVSTE	TVENL-----SSSE	-----	-----	-----	-----	ESVTQVNK	QKPEGVKHEEQQ-RE
opossum	MKLLILSCLVALAVARP--	-----	MVEKI-----SETE	-----	-----	-----	-----	EFVTVIPE	QQI----RREDVPVK
possum	MKLLILTCLVVLAVARP--	-----	MVEKI-----SESE	-----	-----	-----	-----	EHVTDVPE	-----
platypus 2a	MKVFLISCLLAVALAMAMP--	-----	--KL---QSSSSSSSE	ETDQLLVK	EKLVKRRELM	-DLPTTLSSSE	-----	-----	-EHVMEEKEFYQPR
echidna 2a	MKVFLILACLVAVAMALP--	-----	--KQ---HSSSSSSSE	ESDRLLVK	EKLMRRRKLM	-DIPTAFSSSE	-----	-----	-EHSVDPKELYEPQ
platypus 2b	MKVFLILACLVAAVAVPVST	-----	-----	EFDKLLVK	EKLLKHRLDV	KDLPTIFSSE	-----	-----	WEQFLRHPEVYVPLE
echidna 2b	MKVFLIFACLVAVAMAVP--	-----	--KQ---QSSSSSSSE	ETDKQLVM	ENLLKHRALV	KDIPPTFSSE	ENINYEKQ	WEQLLRQPMVYEPFE	-----
	## # # #								

cow	DELQDKIHPFAQTQSLVYPFP--GPIPN-SLPQNIPPLTQTPVVVP--PFLQPEVMGVSKVKEAMAPKQKEMPFPKYP-VEPFTEQSLSLT-----TDVENLHLPPLP
buffalo	DELQDKIHPFAQTQSLVYPFP--GPIPK-SLPQNIPPLTQTPVVVP--PFLQPEIMGVSKVKEAMAPKHKEMPFPKYP-VEPFTEQSLSLT-----TDVENLHLPPLP
goat	DELQDKIHPFAQAQSLVYPFT--GPIPN-SLPQNILPLTQTPVVVP--PFLQPEIMGVVKETMVPKHKEMPFPKYP-VEPFTEQSLSLT-----TDVEKLHLPPLP
sheep	DELQDKIHPFAQAQSLVYPFT--GPIPN-SLPQNILPLTQTPVVVP--PFLQPEIMGVVKETMVPKHKEMPFPKYP-VEPFTEQSLSLT-----TDVEKLHLPPLP
pig	DELQDKIHPFAQAQSLVYPFT--GPIPN-SLPQNILPLTQTPVVVP--PFLQPEIMGVVKETMVPKHKEMPFPKYP-VEPFTEQSLSLT-----TDVEKLHLPPLP
camel	NERQNKIHQFPQPQPLAHPYT--EPIPYPILPQNILPLAQPVFPVV--PLLHPEVMKDSAKETIVPKRKGMPPFKSP-AEPFVEGQSLSLT-----TDFEVLSLPL-
horse	DEQQDKIYTTFPQPQPSLVSHT--EPIPYPILPQNFLPLQPAVMVV--PFLQPKVMDVKPKTKETIIPKRKEMLLQSP-VVPTESQSLSLT-----TDLENLHLPPLP
dog	VERQDKISRFVQPQPVVYPYA--EVPVYAVVPQPSILPLAQPPLIL---PFLQPEIMEVSQAKEТИLPKRKVMPFLKSP-IVPFSERQILNP---TNGENLRLPVH
human	DERQNKIHPLFQQQPLVSPYA--DPIHYAILPQNILPLAQPAVVV---PFLQPEIMEVPVKKENIFPRHKVMPFLKSP-VTPFLDSQILNV---ADLENVHFLPLP
rabbit	DEHQDKIYPSFQPQPLIYPFV--EPIPYGFLPQNILPLAQPAVVL---PVPQPEIMEVPKAKDVTYTKGRVMPVLKSP-TIPFFDPQIPKLT-----TDLENLHLPPLP
rat	EKLQDKILPFIQS---LFPFA--ERIPYPTLPQNILNLAQDMLL---PILQPEIMEDPKAKETIIPPKHKLMPFLKSPKTVFVDSQILNL---REMKNQHLLLP
mouse	DVLQNKFHSGIQSEPKAIYPY--QTISCSPIPQNIQPIAQPPVV-PTDGPII S PELESFLKAKATVLPKHKQMPFLNSETVRLRFNSQIPSL---DLANLHLPQS
elephant	DVLQAKVHSSIQSQPQAFPYQAQATI S CNPVPQNIQPIAQPPVV-PSLGPV I PELESFLKAKATILPKHKQMPFLNSETVRLRINQSIPSL---ASLANLHLPQS
opossum	DEHQNKIQPLFQPQPLVYPFA-EPIPYTVFPPNAIPLAQPIVVL---PFPQPEVKQLPEAKEITFPRQKLMFLKSP---VMPFFDPQIPNLTDLLENLHLPPLP
possum	NERHPEINRFIPLAETMSF-YVPVYWEEMRDAKMTSPLKEKRMTLANI A PEEEELPHI Q HKSLSLAKQRFLASLRP-KAAQFFYAPRMAPLPHKLFTMPK E Q-
platypus 2a	NEHRLEINRYLRPEYEYMMNLYYQPFYWSSEMRNLKMTSLPKDRRMAVLKS V VSDDMLPPLQHKSLSL P KPKVPLPLSH--RQILPPHTLRMVPLSHKLFTIPKREM-
echidna 2a	-----KYPYPFFPPPIKTYVNPHIYQKPAVL P VTHP E LTYLQPQQNPEDMPLP---KKEVLPY L KAVVV V PQVQVMPY P ETEVMPYFPPMTMSLVQPDIV
	-----SYSYPWQS R PINTTY R Y R AYQ I PAVL P MTHPQ T LTYLQPQFKPEDMSISQ--KQIP Y VQAVV M PQ V EAIP P FGAEFMPY A Q P ITTPLLQPEVF
	#

cow	--LLQSWMHQPHQPLPPTV-MFPPQSVLSSLQSKVLPVPQKAVPY--PQRDMPIQAFLLYQE	V-----I
buffalo	--LLQSWMHQPPQPLPPTV-MFPPQSVLSSLQSKVLPVPQKAVPY--PQRDMPIQAFLLYQE	V-----I
goat	--LVQSWMHQPPQPLSPTV-MFPPQSVLSSLQPKVLPVPQKAVPV--QRDMPIQAFLLYQE	V-----LV
sheep	--LVQSWMHQPPQPLPPTV-MFPPQSVLSSLQPKVLPVPQKAVPV--QRDMPIQAFLLYQE	V-----L
pig	--L-QSLMHQIPQPVPQTP-MFAPQPPLSLPQAKVLPVPQQVVPF--PQRDMPFQALLLYQD	V-----
camel	--LQSLMYQIPQPVPQTP-MIPPQSLLSLSQFKVLPVPQQMVPY--PQRAMPVQAVLPFQE	A-----I
horse	--LIQPQFMHQVPQSLLQTL-MLPSQVLSPPQSKVAPFPQPVVVPY--PQRDTPVQAFLLYQD	V-----VI
dog	--LSPLLQPLMHQIPQPPLLQPLMHQIPQPPLPQTP-MLTPQSVLSSIPQPKVLPFPQQVVPY	V-----AL
human	--LLQPLMQQVQPQPIPQTL-ALPPQPLWSVPQPKVLPPIPQQVVPY--PQRAVPVQALLLNQ	V-----
rabbit	--QLLPFHMHQVFQPFQTP-IPYPQALLSLSLPQSKFMPIVPQVVPY--PQRDMPIQALQLFQE	V-----
rat	--PAQLQAQIVQAFP-QTPAVVSSQQLSHPSKSQYLVQQLAPL-FQQGMPVQDLLQYLDLLLNP	V-----
mouse	--LVQLLAQVQAFP-QTH-LVSSQTQLSLPQSKVLYFLQQVAPF--LPQDMSVQDLLQYLE-LLN	V-----
elephant	--LLQPLRHLHQPLAQTP-VLP---LPLSLPKVLPVPQQVVPY--PQRGRPIQNLQLYEEPLLD	V-----
opossum	--LPIAKRDMLSAAELVIPAVH-----ERVIPAIDKREPLPLLAREMPA---LPDKE-----	K----Y
possum	--LPISERERLPA-HERENLLAHEREILLAP-QREMSLIPEREILLAERVVLP---EQEREIRPD	K----N
platypus 2a	--PPSFYREAVIQQLAVPFVR-RESAL--PHQRAIVPVATAAAAV--RESLPLVQQ--EVVPPIMPLDV	KIPETN
echidna 2a	--SAPFYREAV-----LFQERV--LPLHRREIVPPYQRDTIA--RREILPVDQ--RELMPEVVAVDLYPFFQ	KIPETD

Data S6. Alignment of amino acid sequences of CSN3 genes

platypus	-MKTLLLGVAILAMTVGFS	VAEEQKWKRLD	SSESEERWWRLRLKPSLLFRVQDKP-----ERNIPRPSYPYPLLNVPHPNAINPEHQRPYVLP--RFNF-QIPN
echidna	-MKTLLLGVGILVMTVCFS	AAEDEEWKKVD	YSESEERWLRLRKQPSFPFSFGK-----ERNIPRPPYPRPFLNIIPRYTINPEHQFAYVFPNLKF---QIPS
possum	-MKVLFLTWHILAVMVCFS	TADL-DWEKWP	CDKQNERQ-SELQQ-----PLRSPVQYVYTPYHQ-SYVPVIYPRAYVRHPYFSRVAWQKPYPYSYMLPS
opossum	-MKVLFLLIGHILLAMVCFS	TAEL-DWRKWP	CEKQMERP-SELEQQ-----PPGQPPVQDVYTRYTRQ-IYVPILYAPKTSIQYPYFSKLAWQRPYAAVYIPLSS
pig	MMKSSFLIVPILALTLPLFL	GAEEQNQEQLT	RC-ESDKRLFNEEKVKYIPIYYMLNRFPSYGF-FYQRSAVSPNRQFIPYPPYARPVVAGPHAKQPWQ---DQPN
human	-MKSFLVVNNALALTLPLFL	AVEVQNQKOPA	CH-ENDERPFYQKTAAPYVPMYYVPNSYPYGTNLYQRRPAIAINNPYVPRYYANPAVV--RPHAQIPQRQYLPN
rat	MMRNFIIVVMNILALTLPLFL	AAEVQNPDSDN-	CR-EKNEVVYDQVRVLYTPTVSSVLNRN-HYEPIYYHYRTSVP--VSPYAYFPVGLKLLLL-RSPAQILKWQMPN
mouse	MMRNFIIVVNILALTLPLFL	AAEIQNPDSDN-	CRGEKNDIVYDEQRVLYTPVRSVLNFN-QYEPNYHYRPSLPTASPYMYPPLVRLLLL-RSPAPISKWQSMVN
dog	MMKRFFLVVNNIVALALPFL	GAEVQNQEOPT	CR-ENDERLFNQKTVKYIPIHYVLNSFSHYEPNYPHRPAEPINH-Q-YVPYPFYAKPAVAVRTHAQIPQWQVLPN
rabbit	MMKHFLVVNNILAVTLPLFL	AADIQNQEQT	CR-ENEERLFHQVTAPYIIPVHYVMNRYQYEPSPYLLRQAVPTLN-P-FMLNPYYYVKPIV-FKPNVQVPHWQILPN
horse	-MKSFLVVNNILALTLPLFL	GAEVQNQEOPT	CH-KNDERFFDLKTVKYIPIYYVLNSSPRYEPIYYQHLLALLINNQ-HMPYQYYARPA-VRPHVQIPQWQVLPN
elephant	MMKGFLVVNNILLPLPFL	AAEVQNQEESR	CL-EKDERWFCKAVKYIIPNDYVLKSYYRYEPYNQFRAAVPINP-YLTYLYPAKQVA-VRPHTQIPQWQVPSN
camel	-MKSFFLVTILALTLPLFL	GAEVQNQEOPT	CF-EKVERLLNEKTVKYFPIQFVQSRYPSYGINYYQHRLAVPINQ-FIIPYPNYAKPVA-IRLHAQIPQCQALPN
cow	MMKSFFLVTILALTLPLFL	GAQEQNQEQUI	R-CEKDERFFSDKIACYIPIQYVLSRYPSYGLNYYQQRPVALINNQ-FLPYPPYYAKPAA-VRSPAQILQWQVLPN
sheep	MMKSFFLVTILALTLPLFL	GAQEQNQEQUI	C-CEKDERFFDDKIACYIPIQYVLSRYPSYGLNYYQQRPVALINNQ-FLPYPPYYAKPVA-VRSPAQTLQWQVLPN
goat	MMKSFFLVTILALTLPLFL	GAQEQNQEQUI	C-CEKDERFFDDKIACYIPIQYVLSRYPSYGLNYYQQRPVALINNQ-FLPYPPYYAKPAA-VRSPAQILQWQVLPN
buffalo	MMKSFFLVTILALTLPLFL	GAQEQNQEQUI	R-CEKEERFFNDKIACYIPIQYVLSRYPSYGLNYYQQRPVALINNQ-FLPYPPYYAKPAA-VRSPAQILQWQVLPN
	#	#	

platypus	IIP-----FLMFPE--LPPPFFPIVHPIYYDPQTPTTPR-----NPPVTSQTPQPPVDSANT-PEPPTTAPLTATPEAQTPL-QP
echidna	V-----FPFPE--FLPPFYFPVHPIYYGPQTSTPPR-----NPTVTSQTPQPPVHSANT-PESATAAPVTATPEAQTPL-QP
possum	IYP-----WSVVSRLNLHPAFAFNPPHYAQLPVPSSPTNSPTTTIQTNTNIPITNPTSTIVTPAVSSKSAATEDSAAAAMLTSPAAQMA-----
opossum	RYP-----WPVIPR-SPHPSFAFNPPQYARVVPAPSPTSSPAAPMETTTIPTST-STVAATVTPDATSKFVITTEYSTATIPTSPIEQ---QP
pig	VYP-----PTVARRP-RPHASFIAIPPKKNQDKTA-IPAINTSIATVEPTIPAT--EPIVNAEPIVNAVVTPEA\$SEFLITSAPETTVQVTSPVV
human	SHP-----PTVVRPRLHPSFIAPPKKIQDKII-IPTIINTIATVEPTPAPAT--EPTVDSVVTPEAFSESIIITSTPETTTVAVT-----PPTA
rat	FPQ-----PVGVPHPIPNPSFLAIPTNEKHDNTA-IPASNNTIAPIVSTPVSTT-ESVNTVANTEAST--VPISTPETATVPVT-----SPAA
mouse	FPQ-----SAGVPYAIPNPSFLAMPTNENQDNTA-IPTIIDPITPIVSTPVPM-ESIVNTVANPEAST--VSINTPETTTVPVS-----STAA
dog	AYP-----PTMMHRLPQHPSFIAPPKKIQDKTS-IPTIINTIATAEATPIEL--EPKVNTAVTSDASSEFTITSTPETTTVPT-----SPVV
rabbit	IHQ-----PKVGRHSHPFMAILPNKMQDKAV-TPTTNTIAAVEPTPIPTT-EPVV\$TEVIAEASPELII--\$PETTTEATAA----SAAA
horse	IYP-----STVVRHPCPHPSFIAPPKKLQEITV-IPKINTIATVEPTPIPTP-EPTVNNNAVIPDASSEFIIASTPETTTVPTSPV--VQKL
elephant	IYP-----SPSPVPHTYLKPPFIVIPPKKTQDKPI-IPPTGTVASIEATV--EPKVNTVVNAEASSEFIATNTPEATTVPVI-----SPQI
camel	IDP-----PTVERRP-RPRPSFIAPPKKTQDKTV-NPAINTVATVEPPVPITA--EPAVNTVVIAEASSEFITTSTPETT-VQIT-----STEI
cow	TVPAKSCQAQPTTMARHPHPHLSFMAIPPKKNQDKTE-IPTIINTIASGEPTSTPTT-EAVESTVATLED\$PEVIESP-PEINTVQVT-----STAV
sheep	AVPAKSCQDQPTAMARHPHPHLSFMAIPPKKDQDKTE-IPAIINTIASAEPTVHSTPTTEAVVNAVDNPEAESSESIASAP-ETNTAQVT-----STEV
goat	TVPAKSCQDQPTTLARHPHPHLSFMAIPPKKDQDKTE-VPAINTIASAEPTVHSTPTTEAIVNTVDNPEAESSESIASAS-ETNTAQVT-----STEV
buffalo	TVPAKSCQAQPTTMTRHPHPHLSFMAIPPKKNQDKTE-IPTIINTIV\$VEPTSTPTT-EAIENTVATLEASSEVIESV-PETNTAQVT-----STVV
	#

Kawasaki, K., Lafont, A.-G., & Sire, J.-Y. (2011). The evolution of casein genes from tooth genes before the origin of mammals.
Molecular Biology and Evolution, 28, 2053-2061.